

# Supporting Information for “Can we distinguish modes of selective interactions using linkage disequilibrium?”

Aaron P. Ragsdale\*

March 25, 2021

## The diffusion equation and moment system for the two-locus sampling distribution

The two-locus diffusion equation with additive selection was first described by Kimura (1955) and studied extensively in the 1960s and 70s (e.g., Hill and Robertson 1966; Ohta and Kimura 1969). The continuous distribution  $\psi(x_1, x_2, x_3)$  of haplotype frequencies in a population, where  $x_1$  is the frequency of  $AB$ ,  $x_2$  of  $Ab$ , and  $x_3$  of  $aB$ , is governed by the multi-dimensional Fokker-Planck equation:

$$\begin{aligned} \frac{\partial \psi}{\partial \tau} = & \frac{1}{2} \sum_{1 \leq i, j \leq 3} \frac{\partial^2}{\partial x_i \partial x_j} \left[ \frac{x_i(\delta_{i=j} - x_j)\psi}{\nu(\tau)} \right] \\ & - \frac{\rho}{2} \left( -\frac{\partial}{\partial x_1} D\psi + \frac{\partial}{\partial x_2} D\psi + \frac{\partial}{\partial x_3} D\psi \right) \\ & - \frac{\gamma_A}{2} \left[ \frac{\partial}{\partial x_1} x_1(1 - x_1 - x_2)\psi + \frac{\partial}{\partial x_2} x_2(1 - x_1 - x_2)\psi - \frac{\partial}{\partial x_3} x_3(x_1 + x_2)\psi \right] \\ & - \frac{\gamma_B}{2} \left[ \frac{\partial}{\partial x_1} x_1(1 - x_1 - x_3)\psi - \frac{\partial}{\partial x_2} x_2(x_1 + x_3)\psi + \frac{\partial}{\partial x_3} x_3(1 - x_1 - x_3)\psi \right]. \end{aligned} \quad (\text{S1})$$

$D$  is the standard covariance measure of linkage disequilibrium,

$$D = x_1 - (x_1 + x_2)(x_1 + x_3) = x_1 x_4 - x_2 x_3,$$

$\gamma_A$  and  $\gamma_B$  are the (additive) scaled selection coefficients at the left and right locus, and  $\rho$  is the scaled recombination rate between the two loci. Time  $\tau$  is measured in  $2N_e$  generations, and  $\nu(\tau)$  is the population size relative to the ancestral size or some reference size at time  $\tau$ .

Given a function  $\psi$  that solves Equation S1, the two-locus sampling distribution for a sample size of  $n$  haploids can be found by integrating  $\Psi$  against the multinomial sampling function, so that

$$\Psi_n(i, j, k) = \binom{n}{i, j, k, n-i-j-k} \int \int \int_{\substack{x_1, x_2, x_3 \geq 0 \\ x_1 + x_2 + x_3 \leq 1}} \psi(x_1, x_2, x_3) x_1^i x_2^j x_3^k (1 - x_1 - x_2 - x_3)^{n-i-j-k} dx_1 dx_2 dx_3. \quad (\text{S2})$$

In the method-of-moments approach, instead of solving the differential equation for  $\psi$ , we instead integrate both sides of the differential equation against the multinomial sampling function for a given sampling configuration  $(i, j, k)$ . On the left side, we get  $\partial_t \Psi_n(i, j, k)$ , and on the right we obtain, after some simple

---

\*aaronpeaceragsdale@gmail.com

integration by parts and somewhat tedious simplification, terms for drift, recombination, and selection that can be written as sparse linear operators of  $\Psi_n$ . Written compactly, this takes the form

$$\partial_\tau \Psi_n = \frac{1}{2\nu(\tau)} \mathcal{D}_n \Psi_n + \frac{\rho}{2} \mathcal{R}_n \Psi_n + \frac{\theta}{2} \mathcal{U}_n \Psi_n + \mathcal{S}_{n,\gamma,\mathbf{h}} \Psi_n. \quad (\text{S3})$$

Alternatively, we arrive at this same linear system of equations by tracking the expected sampling distribution over  $n$  lineages within the full population and how that changes over time by drawing lineages from one generation to the next in the style of Wright (1931). Both Jouganous et al. (2017), for the single-locus SFS, and Ragsdale and Gravel (2019) drew this connection in detail, so I refer readers to those previous papers for a fuller description of those derivations and discussion. In the next section I simply repeat the results for  $\mathcal{D}$ ,  $\mathcal{R}$ , and  $\mathcal{U}$ , briefly describe the moment closure approximation (which is the same as presented in Ragsdale and Gravel (2019)), and then describe the selection operator  $\mathcal{S}$  for selection with epistasis, both with and without dominance.

## Drift, mutation, recombination, and moment closure

### Drift

Drift for an entry  $(i, j, k)$  depends only on  $\Psi_n$  and therefore closes. The entries of  $\mathcal{D}$  are found by considering the possibility of a coalescence event occurring within a given generation within  $n$  lineages in the full population. If  $n \ll N$ , we can safely assume that at most a single such event occurs in any given generation.

$$\begin{aligned} \mathcal{D}_n(i, j, k) \Psi_n = & (i-1)(n-i-j-k+1) \Psi_n(i-1, j, k) \\ & + (i+1)(n-i-j-k-1) \Psi_n(i+1, j, k) \\ & + (i-1)(k+1) \Psi_n(i-1, j, k+1) \\ & + (i+1)(k-1) \Psi_n(i+1, j, k-1) \\ & + (i-1)(j+1) \Psi_n(i-1, j+1, k) \\ & + (i+1)(j-1) \Psi_n(i+1, j-1, k) \\ & + (j-1)(n-i-j-k+1) \Psi_n(i, j-1, k) \\ & + (j+1)(n-i-j-k-1) \Psi_n(i, j+1, k) \\ & + (j-1)(k+1) \Psi_n(i, j-1, k+1) \\ & + (j+1)(k-1) \Psi_n(i, j+1, k-1) \\ & + (k-1)(n-i-j-k+1) \Psi_n(i, j, k-1) \\ & + (k+1)(n-i-j-k-1) \Psi_n(i, j, k+1) \\ & - 2(i(n-i-j-k) + ik + ij + j(n-i-j-k) + jk + k(n-i-j-k)) \Psi_n(i, j, k) \end{aligned} \quad (\text{S4})$$

### Recombination

If a lineage in our sample of size  $n$  recombines in a given generation, which occurs with probability  $nr$ , we need to draw an extra lineage from the full population for it to recombine with. This means we need  $\Psi_{n+1}$  in the previous generation. After drawing that extra lineage,  $\Psi_n$  changes as we draw one of the two recombinant types (each with probability 1/2) instead of the lineage that was chosen to recombine.

$$\begin{aligned}
\mathcal{R}_n(i, j, k)\Psi_n = & \frac{(i+1)(n-i-j-k+1)}{n+1}\Psi_{n+1}(i+1, j-1, k) \\
& + \frac{(i+1)(n-i-j-k+1)}{n+1}\Psi_{n+1}(i+1, j, k-1) \\
& + \frac{(j+1)(k+1)}{n+1}\Psi_{n+1}(i-1, j+1, k+1) \\
& + \frac{(j+1)(k+1)}{n+1}\Psi_{n+1}(i, j+1, k+1) \\
& - \frac{(i+1)(n-i-j-k)}{n+1}\Psi_{n+1}(i+1, j, k) \\
& - \frac{(j+1)k}{n+1}\Psi_{n+1}(i, j+1, k) \\
& - \frac{j(k+1)}{n+1}\Psi_{n+1}(i, j, k+1) \\
& - \frac{i(n-i-j-k+1)}{n+1}\Psi_{n+1}(i, j, k)
\end{aligned} \tag{S5}$$

### Mutation

We assume an infinite sites mutation (ISM) model where new mutations occur at previously unmutated loci. In the two-locus ISM model, two-locus pairs of variable loci arise when a mutation occurs at one locus when the other locus is already variable. Thus, new mutations at the  $B/b$  locus occur against the single-locus allele frequency distribution  $\Phi_{n,A}$ , and new mutations at the  $A/a$  locus occur against  $\Phi_{n,B}$ , which are found via the single-locus system from Jouganous et al. (2017).

$$\begin{aligned}
\mathcal{U}_n(i, j, k)\Psi_n = & (j+1)\frac{\theta_B}{2}\Phi_{n,A}(j+1)\delta_{i=1, k=0} \\
& + (n-j)\frac{\theta_B}{2}\Phi_{n,A}(j)\delta_{i=0, k=1} \\
& + (i+1)\frac{\theta_A}{2}\Phi_{n,B}(i+1)\delta_{i=1, j=0} \\
& + (n-i)\frac{\theta_A}{2}\Phi_{n,B}(i)\delta_{i=0, j=1}
\end{aligned} \tag{S6}$$

### Jackknife moment closure approximation

We use a jackknife approximation to write the entries of  $\Psi_{n+1}$  and  $\Psi_{n+2}$  as linear combinations of entries in  $\Psi_n$ . The general strategy is to assume the underlying continuous distribution  $\psi(x_1, x_2, x_3)$  can be approximated locally as a quadratic, and then use entries in  $\Psi_n$  that are close in frequency to a given entry in  $\Psi_{n+l}$  to estimate the coefficients of that quadratic using the multinomial sampling formula. Then this quadratically local approximation to  $\psi$  can be used to compute  $\Psi_{n+l}(i, j, k)$  using Eq. (S2). Readers should refer to section S1.3.5 in the Supporting material for Ragsdale and Gravel (2019) for details.

### Selection

First consider the case of no dominance, so that the haplotypes  $Ab$ ,  $aB$ , and  $AB$  have selection coefficients  $s_{Ab}$ ,  $s_{aB}$ , and  $s_{AB}$ , respectively. Note that the case with  $s_{AB} = s_{Ab} + s_{aB}$  implies no epistasis between the  $A/a$  and  $B/b$  loci. Here, we assume all selection coefficients are negative. In a given generation, a selection event could occur in which a haplotype is rejected (selected against) with probability proportional to its

selection coefficient,  $-s$ . We then draw an extra lineage from the full population to replace that rejected lineage.

For example, the probability that an  $AB$  haplotype is selected against and replaced by an  $Ab$  haplotype is

$$-ns_{AB}\frac{i}{n+1}j + 1n\Psi_{n+1}(i, j+1, k),$$

where the additional  $j+1$  lineage in a sample of size  $n+1$  accounts drawing that extra  $Ab$  haplotype. Taking all such selective events together, for additive selection we get

$$\begin{aligned}\mathcal{S}_n(i, j, k)\Psi_n &= \frac{i+1}{n+1}(-s_{AB}(n-i) + s_{Ab}j + s_{aB}k)\Psi_{n+1}(i+1, j, k) \\ &+ \frac{j+1}{n+1}(s_{AB}i - s_{Ab}(n-j) + s_{aB}k)\Psi_{n+1}(i, j+1, k) \\ &+ \frac{k+1}{n+1}(s_{AB}i + s_{Ab}j - s_{aB}(n-k))\Psi_{n+1}(i, j, k+1) \\ &+ \frac{n-i-j-k+1}{n+1}(s_{AB}i + s_{Ab}j + s_{aB}k)\Psi_{n+1}(i, j, k)\end{aligned}\tag{S7}$$

For a general diploid selection model, the idea is nearly the same, but we need to draw an extra lineage to determine the fitness of a diploid individual. For example, the probability that an  $AB$  haplotype is paired with an additional lineage  $Ab$  and selected against, and then replaced by an  $aB$  haplotype is

$$-ns_{AB/Ab}\frac{i}{n+2}\frac{j+1}{n+1}\frac{k+1}{n}\Psi_{n+2}(i, j+1, k+1).$$

There are now many more possible selective events to consider, but after accounting for all possible diploid pairs and replacements (90 in total) and simplifying, we find

$$\begin{aligned}
S_n(i, j, k) \Psi_n = & \frac{n-i-j-k+2}{n+2} \frac{n-i-j-k+1}{n+1} (s_{AB/ab}i + s_{Ab/ab}j + s_{aB/ab}k) \Psi_{n+2}(i, j, k) \\
& + \frac{i+1}{n+2} \frac{n-i-j-k+1}{n+1} (s_{AB/AB}i + s_{AB/Ab}j + s_{AB/aB}k + s_{Ab/ab}j \\
& \quad + s_{aB/ab}k - s_{AB/ab}(n+j+k)) \Psi_{n+2}(i+1, j, k) \\
& + \frac{i+2}{n+2} \frac{i+1}{n+1} (s_{AB/Ab}j + s_{AB/aB}k + s_{AB/ab}(n-i-j-k) - s_{AB/AB}(n-i)) \Psi_{n+2}(i+2, j, k) \\
& + \frac{i+1}{n+2} \frac{j+1}{n+2} (s_{AB/AB}i + s_{AB/aB}k + s_{AB/ab}(n-i-j-k) + s_{Ab/Ab}j \\
& \quad + s_{Ab/aB}k + s_{Ab/ab}(n-i-j-k) - s_{AB/Ab}(2n-i-j)) \Psi_{n+2}(i+1, j+1, k) \\
& + \frac{i+1}{n+2} \frac{k+1}{n+1} (s_{AB/AB}i + s_{AB/Ab}j + s_{AB/ab}(n-i-j-k) + s_{Ab/aB}j \\
& \quad + s_{aB/aB}k + s_{aB/ab}(n-i-j-k) - s_{AB/aB}(2n-i-k)) \Psi_{n+2}(i+1, j, k+1) \\
& + \frac{j+1}{n+2} \frac{n-i-j-k+1}{n+1} (s_{AB/Ab}i + s_{AB/ab}i + s_{Ab/Ab}j + s_{Ab/aB}k \\
& \quad + s_{aB/ab}k - s_{Ab/ab}(n+i+k)) \Psi_{n+2}(i, j+1, k) \\
& + \frac{j+2}{n+2} \frac{j+1}{n+1} (s_{AB/Ab}i + s_{Ab/aB}k + s_{Ab/ab}(n-i-j-k) - s_{Ab/Ab}(n-j)) \Psi_{n+2}(i, j+2, k) \\
& + \frac{j+1}{n+2} \frac{k+1}{n+1} (s_{AB/Ab}i + s_{AB/aB}i + s_{Ab/Ab}j + s_{Ab/ab}(n-i-j-k) \\
& \quad + s_{aB/aB}k + s_{aB/ab}(n-i-j-k) - s_{Ab/aB}(2n-j-k)) \Psi_{n+2}(i, j+1, k+1) \\
& + \frac{k+1}{n+2} \frac{n-i-j-k+1}{n+1} (s_{AB/aB}i + s_{AB/ab}i + s_{Ab/aB}j + s_{Ab/ab}j \\
& \quad + s_{aB/aB}k - s_{aB/ab}(n+i+j)) \Psi_{n+2}(i, j, k+1) \\
& + \frac{k+2}{n+2} \frac{k+1}{n+1} (s_{AB/aB}i + s_{Ab/aB}j + s_{aB/ab}(n-i-j-k) - s_{aB/aB}(n-k)) \Psi_{n+2}(i, j, k+2).
\end{aligned} \tag{S8}$$

Multiplying through by  $2N_{ref}$  gives us selection operators in terms of  $\gamma$  instead of  $s$ .

## Data analysis

### DFE for missense and LOF variants

Loss-of-function (LOF) variants show a dramatic skew toward low-frequency variants across all human populations (Table S4). Here, using the folded SFS for synonymous, missense, and LOF mutations across all autosomal genes, I inferred DFEs for missense and LOF mutations independently. I considered a few different dominance coefficients to explore the effect of the assumed recessivity of the two classes of mutations.

The standard SFS approach to fitting the DFE involves first inferring a demographic history for the population using putatively neutral variants (here, synonymous mutations), and then fixing that demography and fitting a parameterized function for the distribution of selection coefficients for new mutations for the selected classes. DFE inference also requires an estimate for the total mutation rate of the different mutation classes, as much of the signal for strongly selected mutations comes from observing fewer mutations than expected given a known mutation rate (with the assumption that selection purges some fraction of strongly deleterious mutations which are unseen in the sample). Here, I fit demography and DFEs to the folded SFS from the Mende in Sierra Leone (MSL) using *moments* version 1.1.0 (Jouganous et al. 2017).

I used the mutation model from Karczewski et al. (2020) to estimate the total mutation rate across autosomal genes ( $uL$ , where  $u$  is the per-base mutation rate of a given mutation class, and  $L$  is the total length of the coding genome). These values were (0.1442, 0.3426, 0.0256) for synonymous, missense, and LOF mutations, respectively. Roughly two thirds of new mutations in coding regions are expected to be missense mutations,

while only 5% of new mutations are LOF. I fit a demographic model to the synonymous variants, which included a population expansion in the deeper past and exponential growth in the recent past (Figure S7A). Using the inferred optimal scaled mutation rate,  $\theta = 4N_e uL$ , I estimated  $Ne \approx 12,300$ , and assuming an average generation time of 29 years I converted the inferred genetic units to physical units. The best-fit model had a roughly two-fold expansion 400 thousand years ago, and then exponential growth over the past 20-30 thousand years, with a current effective size of  $\approx 63,000$ .

Under this demographic model, I fit a gamma distribution for the distribution of fitness effects to missense and LOF mutations (Table S5). For each fit, I fixed the scaled mutation rate for each mutation class, so that  $\theta_{mis} = \frac{u_{mis}}{u_{syn}} \hat{\theta}_{syn}$  and  $\theta_{lof} = \frac{u_{lof}}{u_{syn}} \hat{\theta}_{syn}$ , where values of  $u$  were found using the GNOMAD mutation model (Karczewski et al. 2020). I tested three values for the dominance coefficient  $h$ : 0, 0.2 and 0.5. For missense mutations,  $h = 0$  gave a poor fit to the data, and  $h = 0.5$  fit best among the three tested dominance coefficients. For LOF variants,  $h = 0$  also fit poorly, but  $h = 0.2$  and  $h = 0.5$  gave similar likelihoods, highlighting that inferring dominance using the SFS is poorly constrained. Regardless of the dominance coefficient assumed, however, the vast majority of LOF variants were inferred to be strongly deleterious, with only  $\sim 10\%$  of new mutations having selection coefficients on the order  $1/N_e$  or less.

## Multinucleotide mutations and positive LD between linked synonymous variants

Multinucleotide mutations (MNMs) are complex mutational events that result in multiple mutations occurring on the same haplotype background in a single generation. Because MNMs fall on the same haplotype, those mutations will be in positive LD, and LD between those pairs that are very tightly linked will not be broken down all that rapidly. MNMs are expected to occur over relatively short distances, on the order of 10s or 100s of base pairs, making them a likely culprit of the observed positive LD among synonymous mutations at short distances.

MNM events can be easily incorporated into the moment system with a simple adjustment to the mutation operator. Instead of all mutations occurring independently in haplotypes with mutations already segregating at the other locus, some fraction of new mutations could instead occur spontaneously and create a new pair of mutations with initial counts  $n_{AB} = 1$  and  $n_{ab} = n - 1$ .

Here, I fit a simple exponential model for the fraction of new mutations at a given distance that arose through a MNM event, so that  $P(MNM|d) = Ae^{-\lambda d}$ , where  $d$  is the distance separating pairs of mutations. I considered all synonymous mutations within genes in the MSL population and used the same population size history model as inferred in the DFE section above for a demographic control. This left two parameters to be fit,  $A$  and  $\lambda$ , which I fit to the binned decay curve of  $\sigma_d^1$ . I needed to assume an average per-base recombination rate  $r$  across gene regions, and tested a number of values between  $10^{-9}$  and  $2 \times 10^{-8}$ . The optimization was insensitive to the chosen value of  $r$ , because the decay of positive LD occurs rapidly. For any plausible value of  $r$ ,  $\sigma_d^1$  decays to zero well before distances between pairs have scaled recombination rates  $\rho = 4N_e r d \sim 1$ , and expected statistics for  $\rho \ll 1$  vary only negligibly.

In fitting the LD decay of  $\sigma_d^1$ , the best fit parameters were  $A = 0.132$ , and  $\lambda = 0.0103$ . An exponential scaling of 0.01 implies that the vast majority of new mutations *pairs* do not occur via MNMs for distances greater than 200 bp, though a substantial fraction (10 – 15%) occur via MNMs for very tightly linked loci with distances on the order 0 – 50 base pairs. It is important to note that this does not mean that 10 – 15% of new mutations occur via MNMs, since this fraction is conditioned on two mutations occurring at short distances.

## Grouping Thousand Genomes populations based on clustering

The large confidence intervals for measurements of signed LD could be driven by either averaging over relatively few observed pairs of mutations, or due to small sample sizes that make each individual measurement a noisy estimate of the LD for that pair of mutations in the full population. To explore the underlying cause of measurement uncertainty in the 1000 Genomes Project Consortium et al. (2015) data, I considered larger sets of samples by combining populations that consistently cluster together in PCA and UMAP space and have low differentiation (Diaz-Papkovich, Patel, and Gravel 2020). I took combinations of CEU/GBR,

CHB/CHS, CDX/KHV, and MSL/GWD. While recognizing that residual population structure in these population combinations could alter expected LD statistics compared to the respective single-population estimates, I was more interested in the effect that increasing the sample sizes would have on estimated measurement error.

Across each of the four combinations tested, confidence intervals were roughly equivalent to those of each of the individual populations. This suggests that the limiting factor to accurate LD measurement is not sample size but rather the overall levels of diversity and number of pairs of mutations that we compare.  $\mathbb{E}[D]$  is most affected by common variants, and the sample sizes of the Thousand Genomes Project data are likely sufficient to accurately estimate common allele frequencies. Adding additional samples will increase the number of rare variants that we observe, but rare variants have minimal impact on  $\sigma_d^1$ . Thus, the accuracy of estimates of  $\sigma_d^1$  is more fundamentally limited by evolutionary history and genome biology (i.e. past population sizes, mutation and recombination rates) than by sample sizes.

## Supplementary Tables

Table S1: General selection model for diploids and dominance models.

Diploid genotype	General model	Simple dominance	Gene-based dominance
$AB / AB$	$1 + s_{AB/AB}$	$1 + 2s_A + 2s_B$	$1 + 2s$
$AB / Ab$	$1 + s_{AB/Ab}$	$1 + 2s_A + 2s_B h_B$	$1 + 2s$
$AB / aB$	$1 + s_{AB/aB}$	$1 + 2s_A h_A + 2s_B$	$1 + 2s$
$AB / ab$	$1 + s_{AB/ab}$	$1 + 2s_A h_A + 2s_B h_B$	$1 + 2sh$
$Ab / Ab$	$1 + s_{Ab/Ab}$	$1 + 2s_A$	$1 + 2s$
$Ab / aB$	$1 + s_{Ab/aB}$	$1 + 2s_A h_A + 2s_B h_B$	$1 + 2s$
$Ab / ab$	$1 + s_{Ab/ab}$	$1 + 2s_A h_A$	$1 + 2sh$
$aB / aB$	$1 + s_{aB/aB}$	$1 + 2s_B$	$1 + 2s$
$aB / ab$	$1 + s_{aB/ab}$	$1 + 2s_B h_B$	$1 + 2sh$
$ab / ab$	1	1	1

Table S2: Haploid epistasis model.

Haplotype	Fitness
$AB$	$(1 + s_A + s_B)(1 + \epsilon)$
$Ab$	$1 + s_A$
$aB$	$1 + s_B$
$ab$	1

Table S3: Thousand Genomes Project population descriptions for populations used in this study.

Code	Description	Region
ESN	Esan in Nigeria	Africa
GWD	Gambian in Western Divisions in the Gambia	Africa
LWK	Luhya in Webuye, Kenya	Africa
MSL	Mende in Sierra Leone	Africa
YRI	Yoruba in Ibadan, Nigeria	Africa
CEU	Utah Residents (CEPH) with Northern and Western European Ancestry	Europe
GBR	British in England and Scotland	Europe
FIN	Finnish in Finland	Europe
IBS	Iberian Population in Spain	Europe
TSI	Toscani in Italia	Europe
CDX	Chinese Dai in Xishuangbanna, China	East Asia
CHB	Han Chinese in Beijing, China	East Asia
CHS	Southern Han Chinese	East Asia
JPT	Japanese in Tokyo, Japan	East Asia
KHV	Kinh in Ho Chi Minh City, Vietnam	East Asia



Table S4: Tamija's  $D$  for classes of coding mutations, both within annotated domains and outside of domains.

Population	Mutation type	Region	Tajima's $D$
ESN	Synonymous	All	-0.882
		In domain	-0.854
		Not in domain	-0.921
	Missense	All	-1.414
		In domain	-1.535
		Not in domain	-1.293
	Loss of function	All	-1.483
		In domain	-2.156
		Not in domain	-1.282
GWD	Synonymous	All	-1.011
		In domain	-0.981
		Not in domain	-1.052
	Missense	All	-1.566
		In domain	-1.678
		Not in domain	-1.452
	Loss of function	All	-1.697
		In domain	-2.328
		Not in domain	-1.501
LWK	Synonymous	All	-1.109
		In domain	-1.088
		Not in domain	-1.139
	Missense	All	-1.589
		In domain	-1.700
		Not in domain	-1.477
	Loss of function	All	-1.666
		In domain	-2.278
		Not in domain	-1.477
MSL	Synonymous	All	-0.983
		In domain	-0.959
		Not in domain	-1.017
	Missense	All	-1.501
		In domain	-1.603
		Not in domain	-1.400
	Loss of function	All	-1.559
		In domain	-2.303
		Not in domain	-1.332
YRI	Synonymous	All	-0.928
		In domain	-0.898
		Not in domain	-0.971
	Missense	All	-1.467
		In domain	-1.586
		Not in domain	-1.348
	Loss of function	All	-1.624
		In domain	-2.237
		Not in domain	-1.424
CEU	Synonymous	All	-0.417
		In domain	-0.392

Table S4: Tamija's  $D$  for classes of coding mutations, both within annotated domains and outside of domains. (*continued*)

Population	Mutation type	Region	Tajima's $D$
FIN	Missense	Not in domain	-0.452
		All	-1.248
		In domain	-1.404
	Loss of function	Not in domain	-1.082
		All	-1.501
		In domain	-2.196
	Synonymous	Not in domain	-1.280
		All	-0.058
		In domain	-0.047
	Missense	Not in domain	-0.075
		All	-0.883
		In domain	-1.048
	Loss of function	Not in domain	-0.710
		All	-1.200
		In domain	-2.034
GBR	Synonymous	Not in domain	-0.906
		All	-0.319
		In domain	-0.300
	Missense	Not in domain	-0.345
		All	-1.120
		In domain	-1.276
	Loss of function	Not in domain	-0.954
		All	-1.313
		In domain	-2.178
	Synonymous	Not in domain	-0.997
		All	-0.689
		In domain	-0.664
	Missense	Not in domain	-0.724
		All	-1.424
		In domain	-1.560
IBS	Loss of function	Not in domain	-1.279
		All	-1.636
		In domain	-2.349
	Synonymous	Not in domain	-1.378
		All	-0.650
		In domain	-0.625
	Missense	Not in domain	-0.685
		All	-1.422
		In domain	-1.568
	Loss of function	Not in domain	-1.266
		All	-1.655
		In domain	-2.349
	Synonymous	Not in domain	-1.397
		All	-0.374
		In domain	-0.366
CDX	Missense	Not in domain	-0.385
		All	-1.179
		In domain	-0.366
	Synonymous	All	-0.374

Table S4: Tamija's  $D$  for classes of coding mutations, both within annotated domains and outside of domains. (*continued*)

Population	Mutation type	Region	Tajima's D	
CHB	Loss of function	In domain	-1.323	
		Not in domain	-1.026	
		All	-1.360	
		In domain	-2.194	
		Not in domain	-1.062	
	Synonymous	All	-0.598	
		In domain	-0.593	
		Not in domain	-0.606	
	Missense	All	-1.389	
		In domain	-1.528	
		Not in domain	-1.239	
	Loss of function	All	-1.586	
In domain		-2.344		
Not in domain		-1.298		
CHS		Synonymous	All	-0.544
			In domain	-0.545
	Not in domain		-0.544	
	Missense	All	-1.334	
		In domain	-1.499	
		Not in domain	-1.150	
	Loss of function	All	-1.559	
		In domain	-2.290	
		Not in domain	-1.292	
JPT		Synonymous	All	-0.371
	In domain		-0.368	
	Not in domain		-0.376	
	Missense	All	-1.194	
		In domain	-1.355	
		Not in domain	-1.019	
	Loss of function	All	-1.410	
		In domain	-2.272	
		Not in domain	-1.086	
KHV		Synonymous	All	-0.576
	In domain		-0.562	
	Not in domain		-0.596	
	Missense	All	-1.346	
		In domain	-1.473	
		Not in domain	-1.210	
	Loss of function	All	-1.535	
		In domain	-2.294	
		Not in domain	-1.269	

Table S5: DFEs inferred for missense and loss-of-function variants in MSL for varying values of  $h$ . General patterns are consistent across different chosen values of  $h$ , although  $h = 0$  results in poorer fits for both missense and LOF variants. Columns to the right of the log-likelihood (LL) column show proportions of new mutations with  $|s|$  in each given bin.

Class	$h$	shape	scale	LL	$[0, 10^{-5})$	$[10^{-5}, 10^{-4})$	$[10^{-4}, 10^{-3})$	$[10^{-3}, 10^{-2})$	$[10^{-2}, \infty)$
Missense	0.0	0.093	768505	-678.2	0.260	0.062	0.077	0.096	0.505
	0.2	0.138	6660	-416.7	0.260	0.098	0.134	0.182	0.327
	0.5	0.147	2117	-392.0	0.282	0.114	0.159	0.214	0.231
LOF	0.0	0.132	99999054	-248.3	0.077	0.028	0.037	0.051	0.807
	0.2	0.177	477994	-226.7	0.083	0.042	0.063	0.095	0.717
	0.5	0.188	121419	-224.2	0.092	0.050	0.077	0.119	0.662

## Supplementary Figures

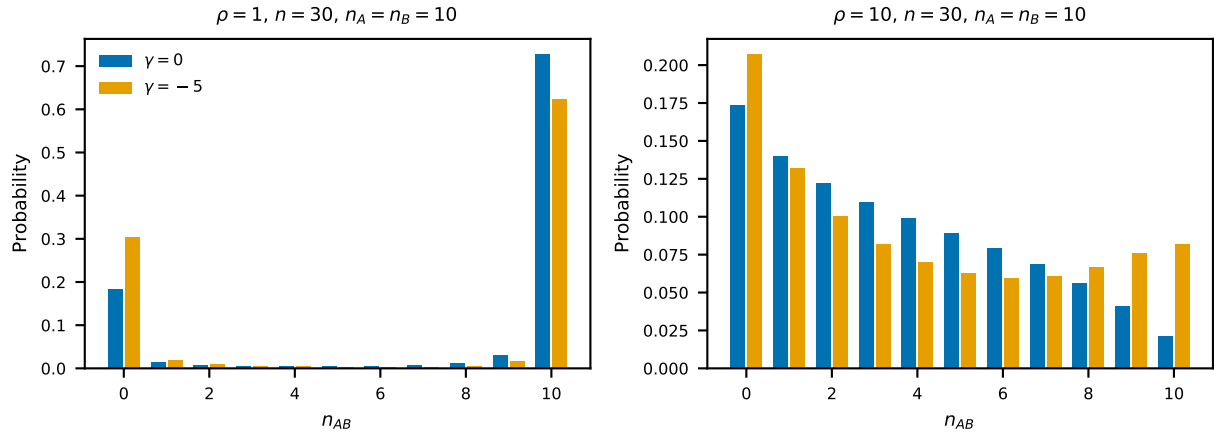


Figure S1: The distribution at stationarity of AB haplotype counts in a sample size of 30, in which we observe 10  $A$  alleles at the left locus, and 10  $B$  alleles at the right locus.

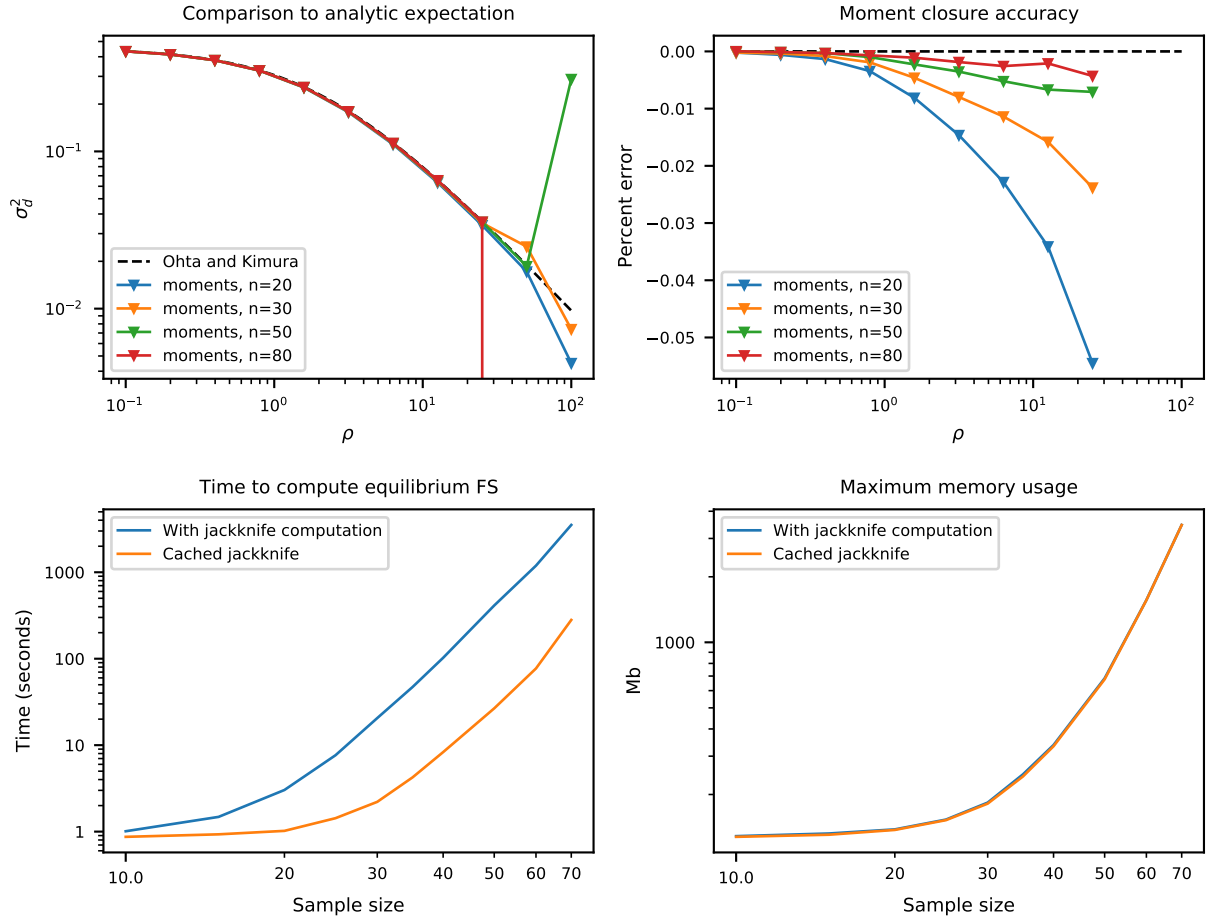


Figure S2: **Accuracy of the jackknife approximation and runtime** Small sample sizes can lead to large error in the closure approximation for larger recombination distances or selection coefficients. Generally, the jackknife approximation breaks down for recombination rates greater than  $\rho \approx 30$ . While increasing sample size leads to more accurate solutions, it comes at the cost of both increased runtime and memory usage. Most analyses performed in this paper used  $n$  between 40 and 70.

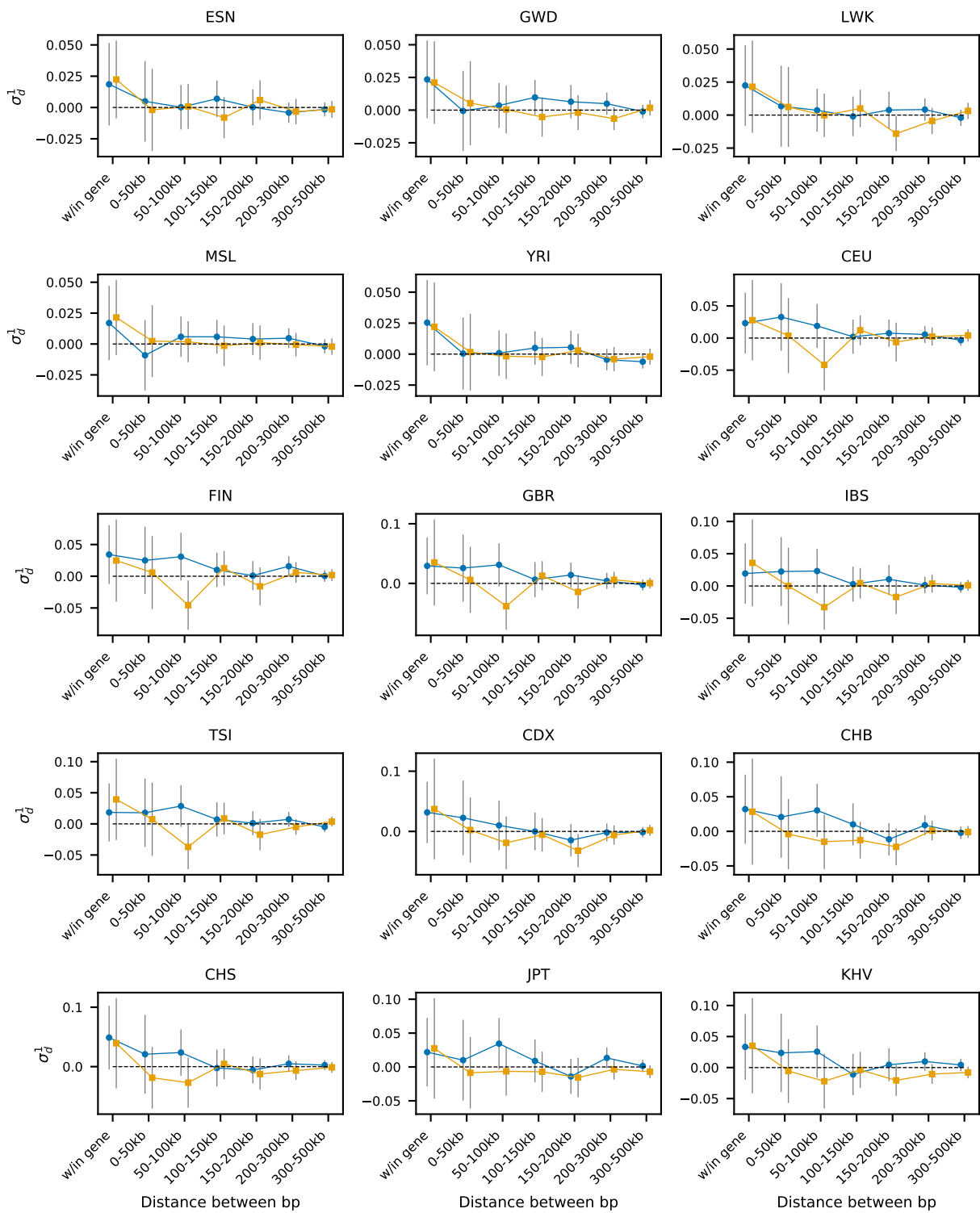


Figure S3: LD within and between protein-coding genes.

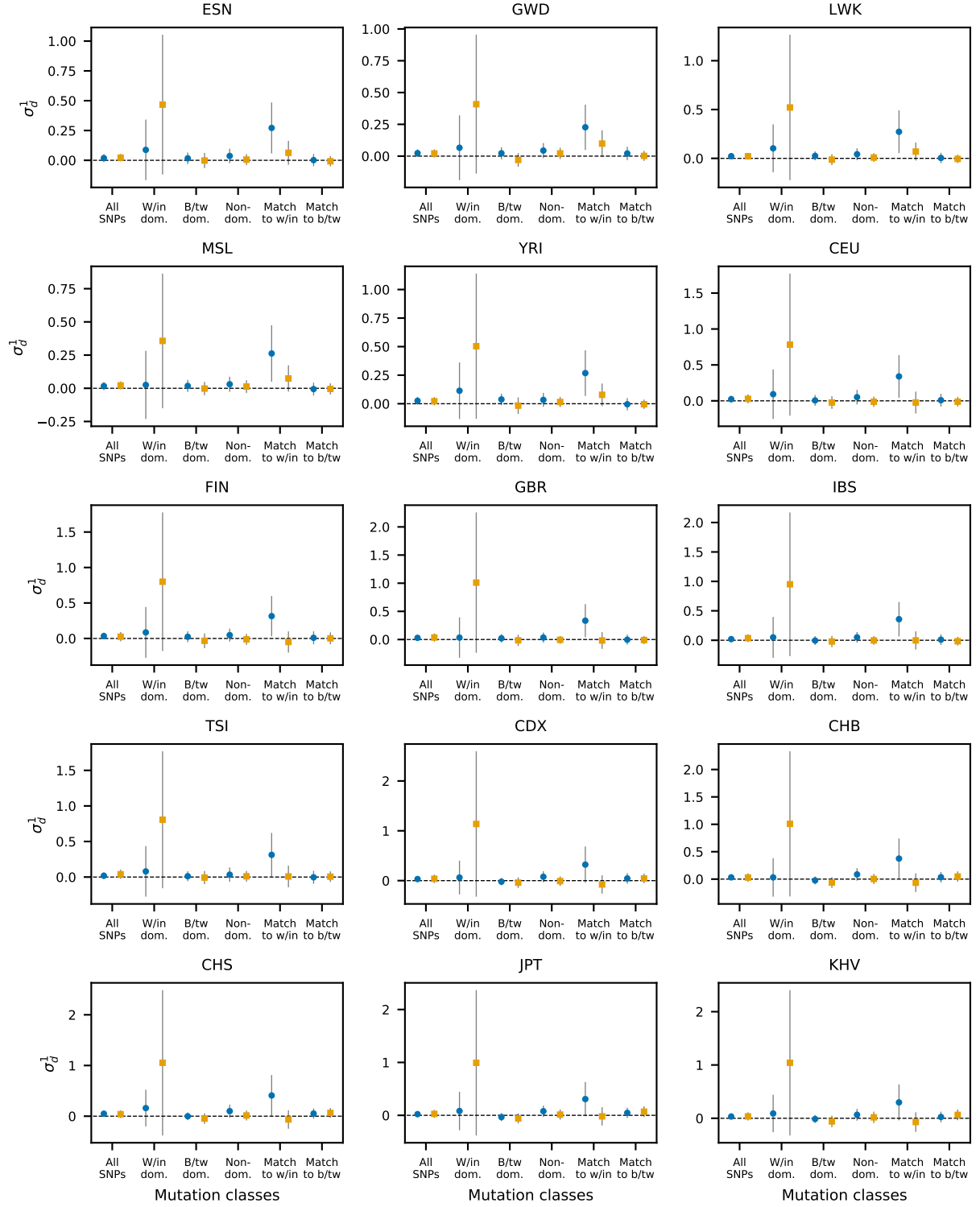


Figure S4: LD within and between coding domains and pairs outside domains at matched distances.



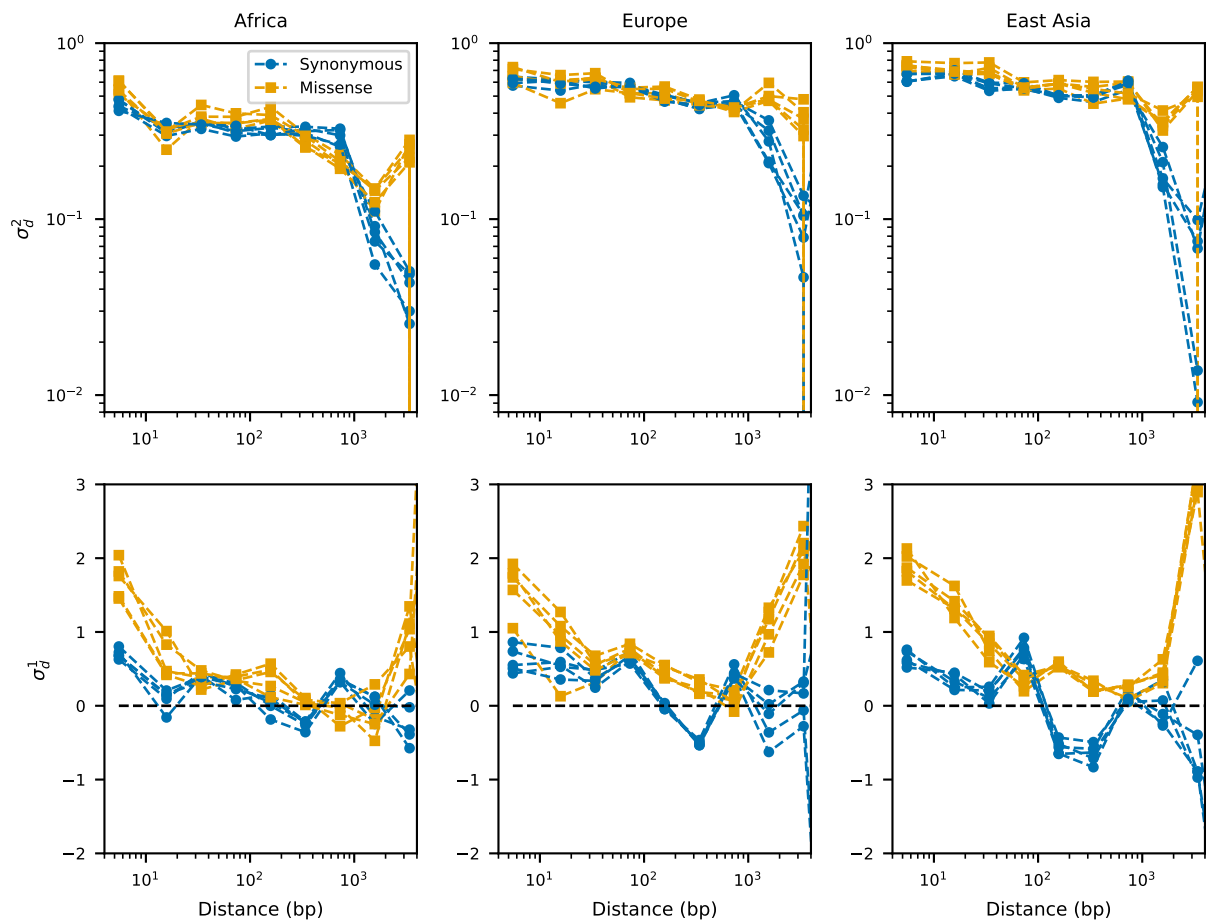


Figure S5: LD decay for synonymous and missense mutations for pairs of mutations that fall inside the same domains.

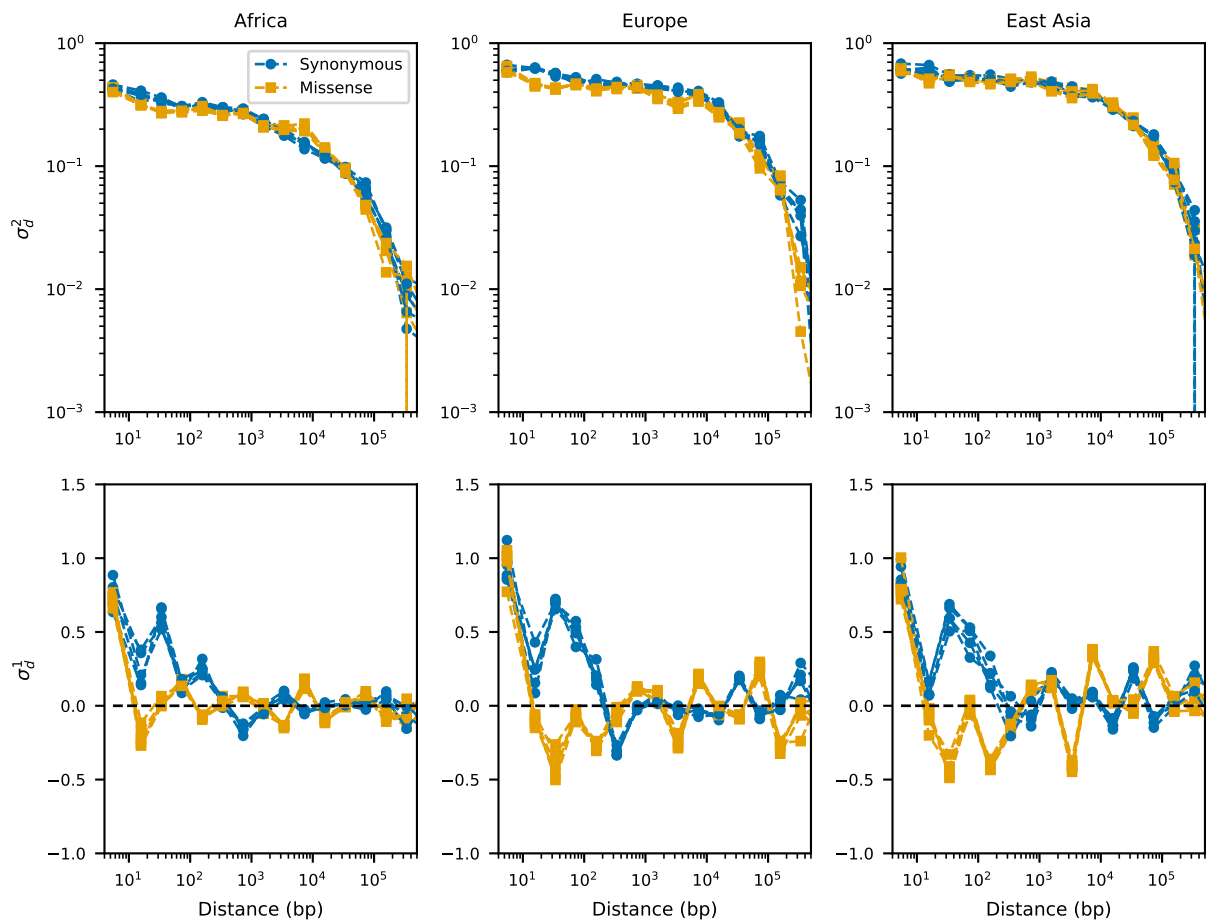


Figure S6: LD decay for synonymous and missense mutations for pairs of mutations that fall outside of domains.

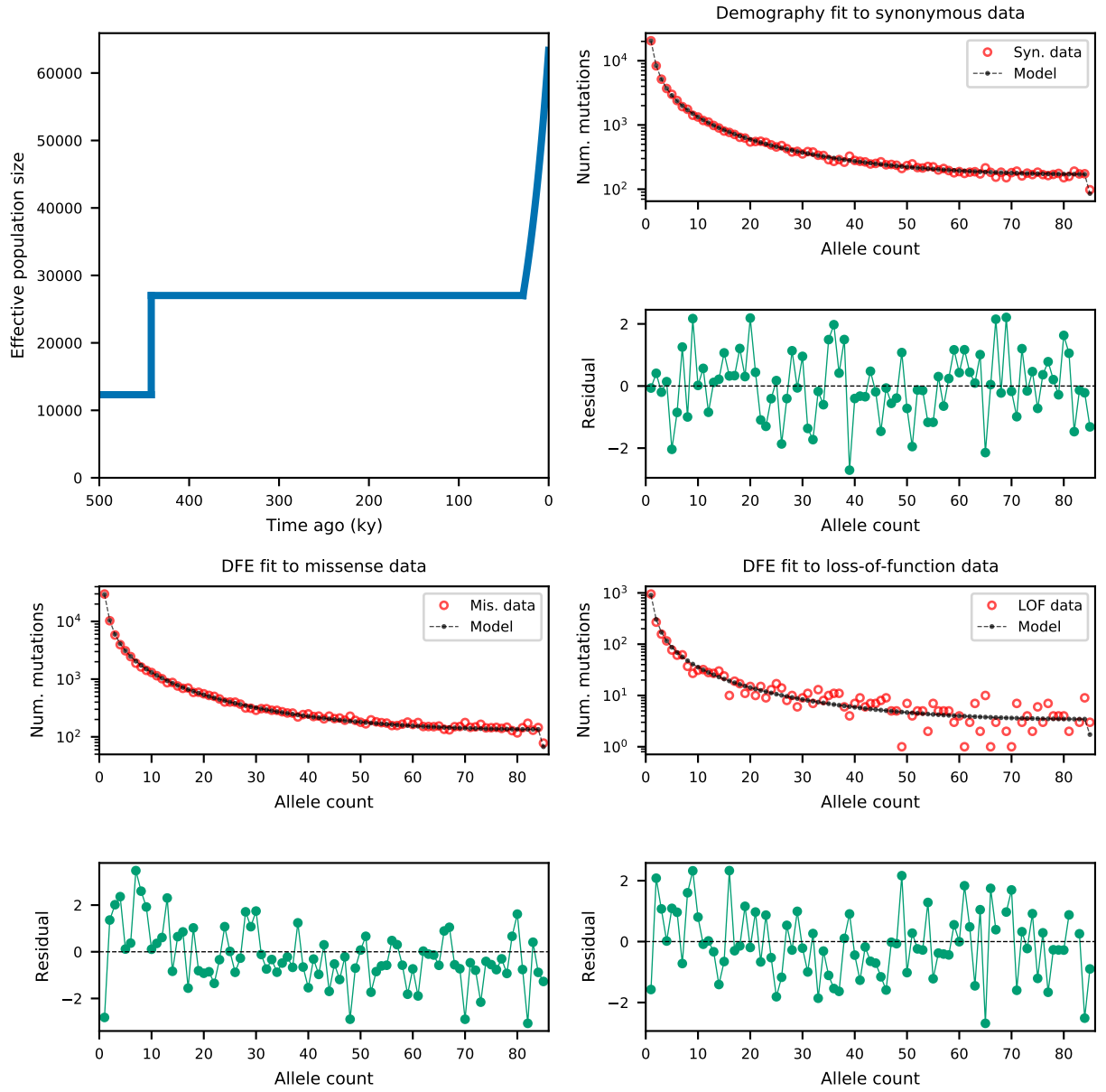


Figure S7: **Demography and DFE for MSL.** A demographic model was fit to the folded synonymous SFS, and DFEs were fit to missense and loss-of-function SFS. Shown here are DFEs fit with  $h = 0.5$ . See Table S5 for inferred best-fit parameters.

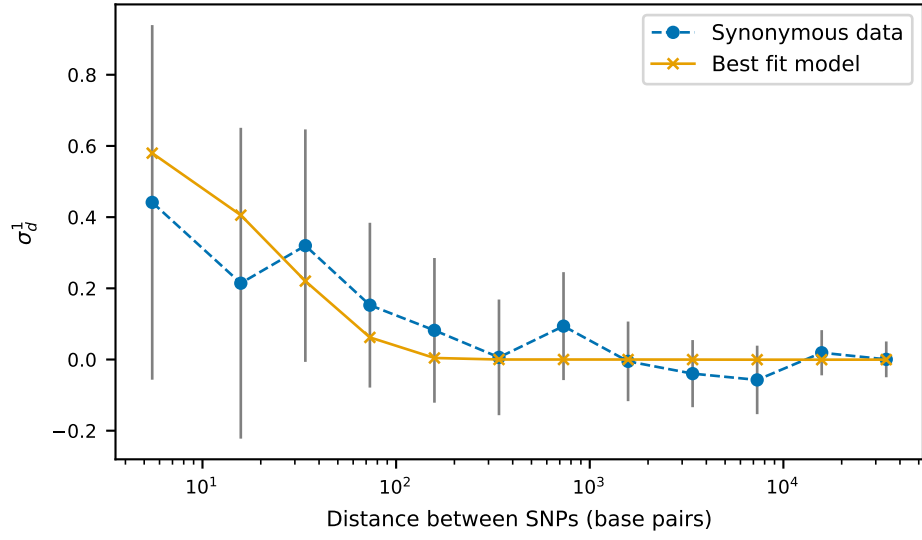


Figure S8: **Optimization of fraction of new mutations arising via multinucleotide mutations by distance.** A simple exponential function was fit to describe the probability that a pair of mutations arose through a MSM event at a given distance  $d$ , as  $Ae^{-\lambda d}$ . Across all recombination rates tested, the best fit parameters were  $A = 0.13$  and  $\lambda = 0.010$ .

## Supporting References

- 1000 Genomes Project Consortium, Adam Auton, Lisa D Brooks, Richard M Durbin, Erik P Garrison, Hyun Min Kang, Jan O Korb, et al. 2015. “A Global Reference for Human Genetic Variation.” *Nature* 526 (7571): 68–74.
- Diaz-Papkovich, Alex, J Patel, and Simon Gravel. 2020. “Density-Based Clustering from Large Multi-Ethnic Biobanks Improves Risk Prediction in Large, Diverse Cohorts.” American Society of Human Genetics conference abstract.
- Hill, W G, and A Robertson. 1966. “The Effect of Linkage on Limits to Artificial Selection.” *Genet. Res.* 8 (3): 269–94.
- Jouganous, Julien, Will Long, Aaron P Ragsdale, and Simon Gravel. 2017. “Inferring the Joint Demographic History of Multiple Populations: Beyond the Diffusion Approximation.” *Genetics* 206 (3): 1549–67.
- Karczewski, Konrad J, Laurent C Francioli, Grace Tiao, Beryl B Cummings, Jessica Alföldi, Qingbo Wang, Ryan L Collins, et al. 2020. “The Mutational Constraint Spectrum Quantified from Variation in 141,456 Humans.” *Nature* 581: 434–43.
- Kimura, Motoo. 1955. “Random Genetic Drift in Multi-Allelic Locus.” *Evolution* 9 (4): 419–35.
- Ohta, T, and M Kimura. 1969. “Linkage Disequilibrium at Steady State Determined by Random Genetic Drift and Recurrent Mutation.” *Genetics* 63 (1): 229–38.
- Ragsdale, Aaron P, and Simon Gravel. 2019. “Models of Archaic Admixture and Recent History from Two-Locus Statistics.” *PLoS Genet.* 15 (6): e1008204.
- Wright, S. 1931. “Evolution in Mendelian Populations.” *Genetics* 16 (2): 97–159.