

# Local PCA Shows How the Effect of Population Structure Differs Along the Genome

Han Li<sup>1</sup>, Peter Ralph<sup>1,2,3,\*</sup>

**1 Department of Molecular and Computational Biology, University of Southern California, Los Angeles, CA , USA**

**2 Institute of Ecology and Evolution, University of Oregon, Eugene, OR, USA**

**3 Department of Mathematics, University of Oregon, Eugene, OR, USA**

\* [plr@uoregon.edu](mailto:plr@uoregon.edu)

## Abstract

Population structure leads to systematic patterns in measures of mean relatedness between individuals in large genomic datasets, which are often discovered and visualized using dimension reduction techniques such as principal component analysis (PCA). Mean relatedness is an average of the relationships across locus-specific genealogical trees, which can be strongly affected on intermediate genomic scales by linked selection and other factors. We show how to use local principal components analysis to describe this meso-scale heterogeneity in patterns of relatedness, and apply the method to genomic data from three species, finding in each that the effect of population structure can vary substantially across only a few megabases. In a global human dataset, localized heterogeneity is likely explained by polymorphic chromosomal inversions. In a range-wide dataset of *Medicago truncatula*, factors that produce heterogeneity are shared between chromosomes, correlate with local gene density, and may be caused by linked selection, such as background selection or local adaptation. In a dataset of primarily African *Drosophila melanogaster*, large-scale heterogeneity across each chromosome arm is explained by known chromosomal inversions thought to be under recent selection, and after removing samples carrying inversions, remaining heterogeneity is correlated with recombination rate and gene density, again suggesting a role for linked selection. The visualization method provides a flexible new way to discover biological drivers of genetic variation, and its application to data highlights the strong effects that linked selection and chromosomal inversions can have on observed patterns of genetic variation.

## 1 Introduction

Wright [68] defined *population structure* to encompass “such matters as numbers, composition by age and sex, and state of subdivision”, where “subdivision” refers to restricted migration between subpopulations. The phrase is also commonly used to refer to the genetic patterns that result from this process, as for instance reduced mean relatedness between individuals from distinct populations. However, it is not necessarily clear what aspects of demography should be included in the concept. For instance, Blair [4] defines *population structure* to be the sum total of “such factors as size of breeding populations, periodic fluctuation of population size, sex ratio, activity range and *differential survival of progeny*” (emphasis added). The definition is similar to Wright’s, but differs in including the effects of natural selection. On closer examination, incorporating differential survival or fecundity makes the concept less clear: should a randomly mating population consisting of two types that are partially reproductively isolated from each other be said to show population structure or not? Whatever the definition, it is clear that due to natural selection, the effects of population structure – the *realized* patterns of genetic relatedness – differ depending on which portion of the genome is being considered. For instance, strongly locally adapted alleles of a gene will be selected against in migrants to different habitats, increasing genetic differentiation between populations near to this gene. Similarly, newly adaptive alleles spread first in local populations. These observations motivate many methods to search for genetic loci under selection, as for example in Huerta-Sánchez et al. [30], Martin et al. [46], and Duforet-Frebourg et al. [19].

These realized patterns of genetic relatedness summarize the shapes of the genealogical trees at each location along the genome. Since these trees vary along the genome, so does relatedness, but averaging over sufficiently many trees we hope to get a stable estimate that doesn’t depend much on the genetic markers chosen. This is not guaranteed: for instance, relatedness on sex chromosomes is expected to differ from the autosomes; and positive or negative selection on particular loci can dramatically distort shapes of nearby genealogies [3, 10, 36]. Indeed, many species show chromosome-scale variation in diversity and divergence (e.g., Langley et al. [40]); species phylogenies can differ along the genome due to incomplete lineage sorting, adaptive introgression and/or local adaptation (e.g., Ellegren et al. [21], Nadeau et al. [49], Pease and Hahn [55], Pool [57], Vernot and Akey [65]); and theoretical expectations predict that geographic patterns of relatedness should depend on selection [14].

Patterns in genome-wide relatedness are often summarized by applying principal components analysis (PCA, Patterson et al. [54]) to the genetic covariance matrix, as pioneered by Menozzi et al. [48]. The results of PCA can be related to the genealogical history of the samples, such as time to most recent common ancestor and migration rate between populations [47, 51], and sometimes produce “maps” of population structure that reflect the samples’ geographic origin distorted by rates of gene flow [52].

Modeling such “background” kinship between samples is essential to genome-wide as-

sociation studies (GWAS, Astle and Balding [2], Price et al. [59]), and so understanding variation in kinship along the genome could lead to more generally powerful methods, and may be essential for doing GWAS in species with substantial heterogeneity in realized patterns of mean relatedness along the genome.

Others have applied PCA to windows of the genome: Ma and Amos [43] used local PCA much as we do to identify putative chromosomal inversions. Bryc et al. [7] and Brisbin et al. [6] used PCA to infer tracts of local ancestry in recently admixed populations, but by projecting each genomic window onto the axes of a single, globally-defined PCA rather than doing PCA separately on each window.

A note on nomenclature: In this work we describe variation in patterns of relatedness using local PCA, where “local” refers to proximity along the genome. A number of general methods for dimensionality reduction also use a strategy of “local PCA” (e.g., Kambhatla and Leen [34], Manjón et al. [45], Roweis and Saul [60], Weingessel and Hornik [67]), performing PCA not on the entire dataset but instead on subsets of observations, providing local pictures which are then stitched back together to give a global picture. At first sight, this differs from our method in that we restrict to subsets of *variables* instead of subsets of observations. However, if we flip perspectives and think of each genetic variant as an observation, our method shares common threads, although our method does not subsequently use adjacency along the genome, as we aim to identify similar regions that may be distant.

It is common to describe variation along the genome of simple statistics such as  $F_{ST}$  and to interpret the results in terms of the action of selection (e.g., Ellegren et al. [21], Turner et al. [64]). However, a given pattern (e.g., valleys of  $F_{ST}$ ) can be caused by more than one biological process [8, 18], which in retrospect is unsurprising given that we are using a single statistic to describe a complex process. It is also common to use methods such as PCA to visualize large-scale patterns in mean genome-wide relatedness. In this paper we show if and how patterns of mean relatedness vary systematically along the genome, in a way particularly suited to large samples from geographically distributed populations. Geographic population structure sets the stage by establishing “background” patterns of relatedness; our method then describes how this structure is affected by selection and other factors. Our aim is not to identify outlier loci, but rather to describe larger-scale variation shared by many parts of the genome; correlation of this variation with known genomic features can then be used to uncover its source.

## 2 Materials and Methods

As depicted in Figure 1, the general steps to the method are: (1) divide the genome into windows, (2) summarize the patterns of relatedness in each window, (3) measure dissimilarity in relatedness between each pair of windows, (4) visualize the resulting dissimilarity matrix using multidimensional scaling (MDS), and (5) combine similar windows to more

accurately visualize local effects of population structure using PCA.

## 2.1 PCA in genomic windows

To begin, we first recoded sampled genotypes as numeric matrices in the usual manner, by recording the number of nonreference alleles seen at each locus for each sample. We then divided the genome into contiguous segments (“windows”) and applied principal component analysis (PCA) as described in McVean [47] separately to the submatrices that corresponded to each window. The choice of window length entails a tradeoff between signal and noise, since shorter windows allow better resolution along the genome but provide less precise estimates of relatedness. A method for choosing a window length to balance these considerations is given in Appendix A. Precisely, denote by  $Z$  the  $L \times N$  recoded genotype matrix for a given window ( $L$  is the number of SNPs and  $N$  is the sample size), and by  $\overline{Z}_s$  the mean of non-missing entries for allele  $s$ , so that  $\overline{Z}_s = \frac{1}{n_s} \sum_j Z_{sj}$ , where the sum is over the  $n_s$  nonmissing genotypes. We first compute the mean-centered matrix  $X$ , as  $X_{si} = Z_{si} - \overline{Z}_s$ , and preserving missingness. (This mean-centering makes the result not depend on the choice of reference allele, exactly if there is no missing data, and approximately otherwise.) Next, we find the covariance matrix of  $X$ , denoted  $C$ , as  $C_{ij} = \frac{1}{m_{ij}-1} \sum_s X_{si} X_{sj} - \frac{1}{m_{ij}(m_{ij}-1)} (\sum_s X_{si})(\sum_s X_{sj})$ , where all sums are over the  $m_{ij}$  sites where both sample  $i$  and sample  $j$  have nonmissing genotypes. The principal components are the eigenvectors of  $C$ , normalized to have Euclidean length equal to one, and ordered by magnitude of the eigenvalues.

The top 2–5 principal components are generally good summaries of population structure; for ease of visualization we usually only use the first two (referred to as  $PC1$  and  $PC2$ ), and check that results hold using more. The above procedure can be performed on any subset of the data; for future reference, denote by  $PC1_j$  and  $PC2_j$  the result after applying to all SNPs in the  $j^{\text{th}}$  window. (Note, however, that our measure of dissimilarity between windows does not depend on PC ordering.)

## 2.2 Similarity of patterns of relatedness between windows

We think of the local effects of population structure as being summarized by *relative* position of the samples in the space defined by the top principal components. However, we do not compare patterns of relatedness of different genomic regions by directly comparing the PCs, since rotations or reflections of these imply identical patterns of relatedness. Instead, we compare the low-dimensional approximations of the local covariance matrices obtained using the top  $k$  PCs, which is invariant under reordering of the PCs, reflections, and rotations and yet contains all other information about the PCs. (For results shown here, we use  $k = 2$ .) Furthermore, to remove the effect of artifacts such as mutation rate variation, we also rescale each approximate covariance matrix to be of similar size (precisely, so that the underlying data matrix has trace norm equal to one).

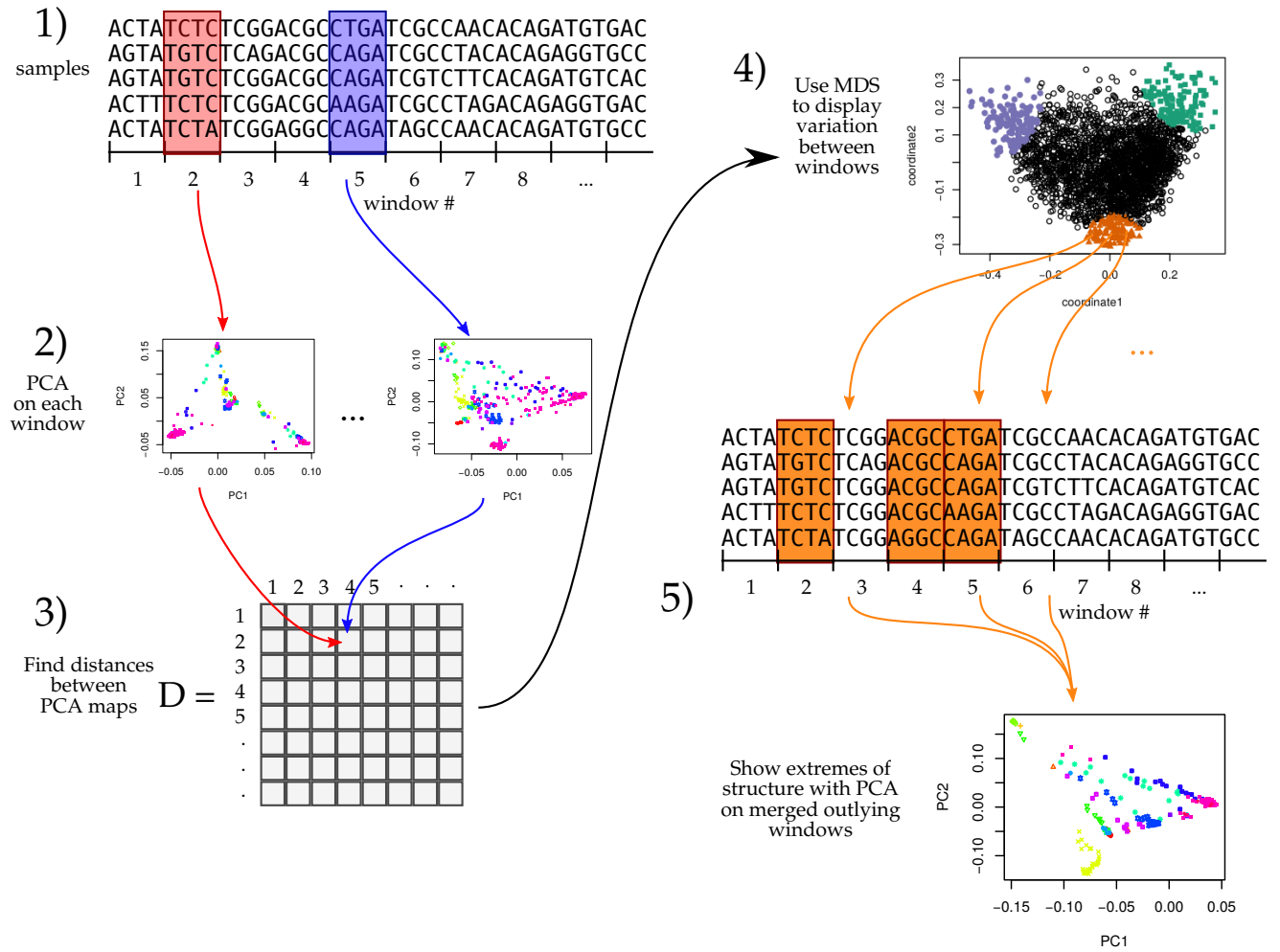


Figure 1: An illustration of the method; see Methods for details.

To do this, define the  $N \times k$  matrix  $V(i)$  so that  $V(i)_{\cdot\ell}$ , the  $\ell^{\text{th}}$  column of  $V(i)$ , is equal to the  $\ell^{\text{th}}$  principal component of the  $i^{\text{th}}$  window, multiplied by  $(\lambda_{\ell i} / \sum_{m=1}^k \lambda_{mi})^{1/2}$ , where  $\lambda_{\ell i}$  is the  $\ell^{\text{th}}$  eigenvalue of the genetic covariance matrix. Then, the rescaled, rank  $k$  approximate covariance matrix for the  $i^{\text{th}}$  window is

$$M(i) = \sum_{\ell=1}^k V(i)_{\cdot\ell} V(i)_{\cdot\ell}^T. \quad (1)$$

To measure the similarity of patterns of relatedness for the  $i^{\text{th}}$  window and  $j^{\text{th}}$  window, we then use Euclidean distance  $D_{ij}$  between the matrices  $M(i)$  and  $M(j)$ , defined by  $D_{ij}^2 = \sum_{k,\ell} (M(i)_{k,\ell} - M(j)_{k,\ell})^2$ .

The goal of comparing PC plots up to rotation and reflection turned out to be equivalent to comparing rank- $k$  approximations to local covariance matrices. This suggests instead directly comparing entire local covariance matrices. However, with thousands of samples and tens of thousands of windows, computing the distance matrix would take months of CPU time, while as defined above,  $D$  can be computed in minutes using the following method. Since for square matrices  $A$  and  $B$ ,  $\sum_{ij} (A_{ij} - B_{ij})^2 = \sum_{ij} (A_{ij}^2 + B_{ij}^2) - 2 \text{tr}(A^T B)$ , then due to the orthogonality of eigenvectors and the cyclic invariance of trace,  $D_{ij}$  can be computed efficiently as

$$D_{ij} = \left( \frac{\sum_{\ell=1}^k \lambda_{\ell i}^2}{(\sum_{\ell=1}^k \lambda_{\ell i})^2} + \frac{\sum_{\ell=1}^k \lambda_{\ell j}^2}{(\sum_{\ell=1}^k \lambda_{\ell j})^2} - 2 \sum_{\ell,m=1}^k (V(i)^T V(j))_{\ell m}^2 \right)^{1/2}. \quad (2)$$

### 2.3 Visualization of results

We use multidimensional scaling (MDS) to visualize relationships between windows as summarized by the dissimilarity matrix  $D$ . MDS produces a set of  $m$  coordinates for each window that give the arrangement in  $m$ -dimensional space that best recapitulates the original distance matrix. For results here, we use  $m = 2$  to produce one- or two-dimensional visualizations of relationships between windows' patterns of relatedness.

We then locate variation in patterns of relatedness along the genome by choosing collections of windows that are nearby in MDS coordinates, and map their positions along the genome. A visualization of the effects of population structure across the entire collection is formed by extracting the corresponding genomic regions and performing PCA on all, aggregated, regions.

### 2.4 Testing

We tested the method using two types of simulation. First, to verify expected behavior, we simulated "genomes" as an independent sequence of correlated Gaussian "genotypes", using a different covariance matrix in the first quarter, middle half, and last quarter of

the chromosome. The details of the simulation, also designed to detect sensitivity to PC switching, are given in Appendix B.1. To verify robustness to missing data, we ran the method after randomly dropping 50% of the genotypes in the first half of the genome; if the method is misled by missing data, then it will distinguish the two halves of the chromosome rather than the segments having different covariance matrices.

To provide a realistic test, we next used forwards-time, individual-based simulations, implemented using SLiM v3 [25], which are described in detail in Appendix B.2. To provide realistic population structure for PCA to identify, each simulation had at least 5,000 diploid individuals, living across a continuous square range, with Gaussian dispersal and local density-dependent competition. Each genome was modeled on human chromosome 7, which is  $1.54 \times 10^8$  bp long, with an overall recombination rate of 1.6785 crossovers per chromosome per generation. To improve speed, we used tskit [35] to record tree sequences in SLiM [26] and to add neutral mutations afterwards, at a rate of  $10^{-9}$  per bp per generation. Most simulations were neutral, but we also included linked selection, of two types. First, we introduced selected mutations into two regions, which extended from 1/3 to 1/2 and from 5/6 to the end of the genome respectively. These had selection coefficients from a Gamma distribution with shape 2 and mean 0.005 at a rate of  $10^{-10}$  per bp, that were either beneficial (with probability 1/30) or deleterious (otherwise). Second, to roughly model a recent expansion followed by local adaptation, we introduced mutations in the same manner as above, except that mutations were no longer unconditionally deleterious or beneficial: each selection coefficient was multiplied by a factor depending on the spatial location of the individual being evaluated, varying linearly from -1 at the left side of the range to +1 at the right edge. In all simulations, genome-wide PCA displayed a map of the population range, as expected.

## 2.5 Datasets

We applied the method to genomic datasets with good geographic sampling: 380 African *Drosophila melanogaster* from the Drosophila Genome Nexus [39], a worldwide dataset of humans, 3,965 humans from several locations worldwide from the POPRES dataset [50], and 263 *Medicago truncatula* from 24 countries around the Mediterranean basin a range-wide dataset of the partially selfing weedy annual plant from the *Medicago truncatula* Hapmap Project [63], as summarized in Table 1.

***Drosophila melanogaster*:** We used whole-genome sequencing data from the Drosophila Genome Nexus (<http://www.johnpool.net/genomes.html>, [39]), consisting of the Drosophila Population Genomics Project phases 1–3 [40, 58], and additional African genomes [39]. After removing 20 genomes with more than 8% missing data, we were left with 380 samples from 16 countries across Africa and Europe. Since the *Drosophila* samples are from inbred lines or haploid embryos, we treat the samples as haploid when recoding; regions with residual heterozygosity were marked as missing in the original dataset; we also removed

positions with more than 20% missing data. Each chromosome arm we investigated (X, 2L, 2R, 3L, and 3R) has 2–3 million SNPs; PCA plots for each arm are shown in Figure S1.

**Human:** We also used genomic data from the entire POPRES dataset [50], which has array-derived genotype information for 447,267 SNPs across the 22 autosomes of 3,965 samples in total: 346 African-Americans, 73 Asians, 3,187 Europeans and 359 Indian Asians. Since these data derive from genotyping arrays, the SNP density is much lower than the other datasets, which are each derived from whole genome sequencing. We excluded the sex chromosomes and the mitochondria. PCA plots for each chromosome, separately, are shown in Figure S2.

***Medicago truncatula*:** Finally, we used whole-genome sequencing data from the *Medicago truncatula* Hapmap Project [63], which has 263 samples from 24 countries, primarily distributed around the Mediterranean basin. Each of the 8 chromosomes has 3–5 million SNPs; PCA plots for these are shown in Figure S3. We did not use the mitochondria or chloroplasts.

species	# SNPs per window	mean window length (bp)	mean # windows per chromosome	mean % variance explained by top 2 PCs
<i>Drosophila melanogaster</i>	1,000	9,019	2,674	0.53
Human	100	636,494	203	0.55
<i>Medicago truncatula</i>	10,000	102,580	467	0.50

Table 1: Descriptive statistics for each dataset used.

## 2.6 Data access

The methods described here are implemented in an open-source R package available at [https://github.com/petrelharp/local\\_pca](https://github.com/petrelharp/local_pca), as well as scripts to perform all analyses from VCF files at various parameter settings.

Datasets are available as follows: human (POPRES) at dbGaP with accession number phs000145.v4.p2, *Medicago* at the Medicago Hapmap <http://www.medicagohapmap.org/>, and *Drosophila* at the Drosophila Genome Nexus, <http://www.johnpool.net/genomes.html>.



### 3 Results

In all three datasets: a worldwide sample of humans, African *Drosophila melanogaster*, and a rangewide sample of *Medicago truncatula*, PCA plots vary along the genome in a systematic way, showing strong chromosome-scale correlations. This implies that variation is due to meaningful heterogeneity in a biological process, since noise due to randomness in choice of local genealogical trees is not expected to show long distance correlations. Below, we discuss the results and likely underlying causes.

#### 3.1 Validation

Simple non-population-based simulations with Gaussian “genotypes” showed that the method performs as expected, clearly separating regions of the genome with different underlying covariance matrices without being affected by extreme differences in amount of missing data (Supplemental Figure S4). This simulation also verifies insensitivity to ordering of top PCs, since it was performed using a covariance matrix with the top two eigenvalues equal, so that the order of empirical eigenvectors (PCs) switches randomly.

Individual-based simulations using SLiM [25] allowed us to test the effects of recombination and mutation rate variation, as well as linked selection. As expected, varying recombination rate stepwise by a factor of 64 did not induce patterns in the MDS visualizations correlated with recombination rate (Supplemental Figure S5). Since varying mutation rate with a fixed recombination map is equivalent to varying the recombination map and remapping windows, this also indicates that the method is not misled by variation in mutation rate. On the other hand, a recombination map with hotspots (the HapMap human female map for chromosome 7 [32]) induced outliers at long regions of low recombination rate (also as expected).

Simulations with linked selection produced mixed results (Figure S6). The method strongly identified the regions under spatially varying linked selection. It also identified the regions (although less unambiguously) with constant selection and stepwise varying recombination rate, but did not clearly identify them with constant recombination rate. This difference is likely because recombination rates are overall lower in the first case, leading to a stronger effect of linked selection. These tests are not meant to be comprehensive survey of linked selection, but only to demonstrate that linked selection can produce signals similar to what we see in real data.

#### 3.2 *Drosophila melanogaster*

We applied the method to windows of average length 9 Kbp, across chromosome arms 2L, 2R, 3L, 3R and X separately. The first column of Figure 2 is a multidimensional scaling (MDS) visualization of the matrix of dissimilarities between genomic windows: in other words, genomic windows that are closer to each other in the MDS plot show more similar patterns of relatedness. For each chromosome arm, the MDS visualization roughly

resembles a triangle, sometimes with additional points. Since the relative position of each window in this plot shows the similarity between windows, this suggests that there are at least three extreme manifestations of population structure typified by windows found in the “corners” of the figure, and that other windows’ patterns of relatedness may be a mixture of those extremes. The next two columns of Figure 2 respectively depict the two MDS coordinates of each window, plotted against the window’s position along the genome, to show how the plot of the first column is laid out along the genome. The patterns did not depend on the number of PCs used (see Figure S7 for the same plot with  $k = 5$  PCs), and are only weakly correlated with variation in missingness (see Figure S8).

To help visualize how clustered windows with similar patterns of relatedness are along each chromosome arm, we selected three “extreme” windows in the MDS plot and the 5% of windows that are closest to it in the MDS coordinates, then highlighted these windows’ positions along the genome, and created PCA plots for the windows, combined. Representative plots are shown for three groups of windows on each chromosome arm in Figure 2 (groups are shown in color), and in Supplemental Figure S9 (PCA plots). The latter plots are quite different, showing that genomic windows in different regions of the MDS plot indeed show quite different patterns of relatedness.

The most striking variation in patterns of relatedness turns out to be explained by several large inversions that are polymorphic in these samples, discussed in Corbett-Detig and Hartl [16] and Langley et al. [40]. To depict this, Figure 3 shows the PCA plots in Supplemental Figure S9 recolored by the orientation of the inversion for each sample. Taking chromosome arm 2L as an example, the two regions of similar, extreme patterns of relatedness shown in green in the first row of Figure 2 lie directly around the breakpoints of the inversion  $\text{In}(2\text{L})\text{t}$ , and the PCA plots in the first rows of Figure 3 shows that patterns of relatedness here are mostly determined by inversion orientation. The regions shown in purple on chromosome 2L lie near the centromere, and have patterns of relatedness reflective of two axes of variation, seen in Figure 3 and Supplemental Figure S9, which correspond roughly to latitude within Africa and to degree of cosmopolitan admixture respectively (see Lack et al. [39] for more about admixture in this sample). The regions shown in orange on chromosome 2L mostly lie inside the inversion, and show patterns of relatedness that are a mixture between the other two, as expected due to recombination within the (long) inversion [24]. Similar results are found in other chromosome arms, albeit complicated by the coexistence of more than one polymorphic inversion; however, each breakpoint visibly affects patterns in the MDS coordinates (see vertical lines in Figure 2).

To see how patterns of relatedness vary in the absence of polymorphic inversions, we performed the same analyses after removing, for each chromosome arm, any samples carrying inversions on that arm. In the result, shown in Figure 4 and Supplemental Figure S10, the striking peaks associated with inversion breakpoints are gone, and previously smaller-scale variation now dominates the MDS visualization. For instance, the majority of the variation along 3L in Figure 2 is on the left end of the arm, dominated by two large peaks

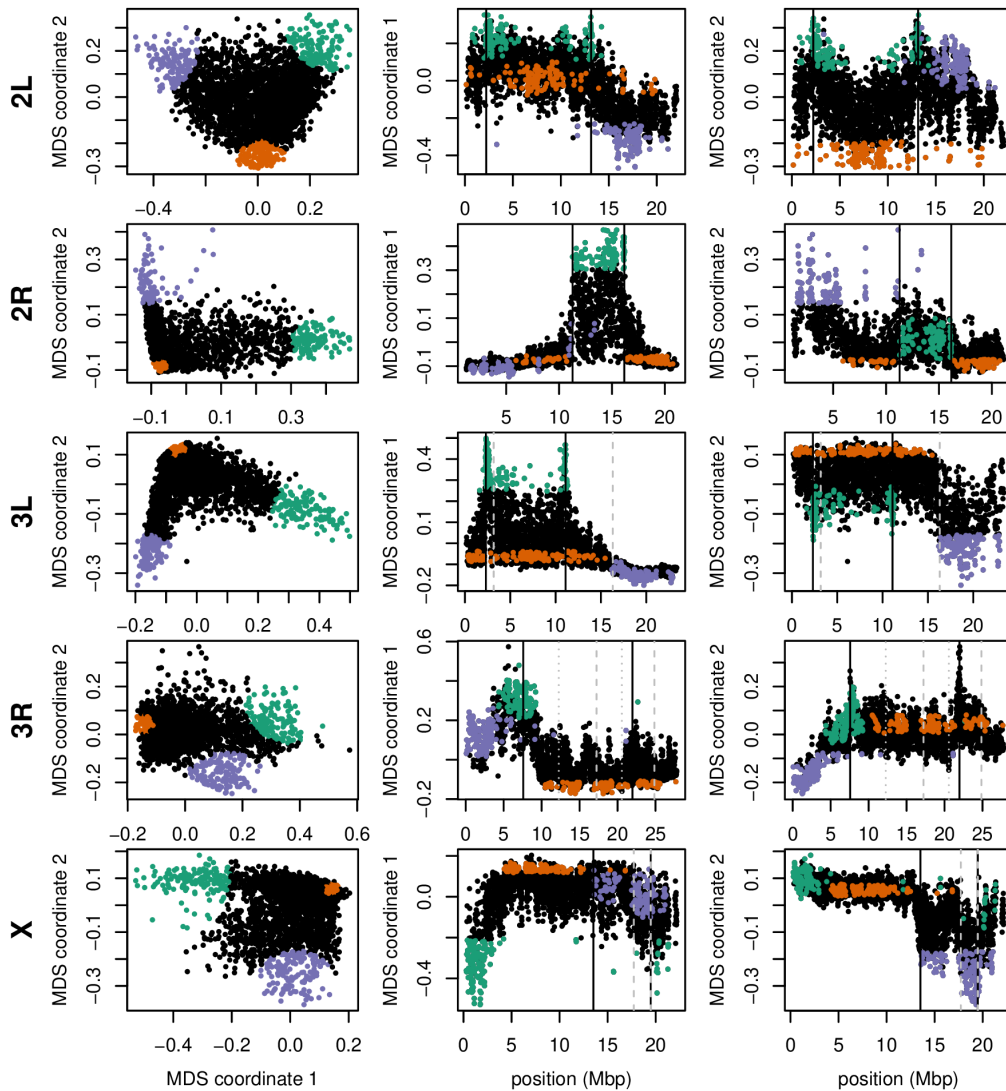


Figure 2: Variation in patterns of relatedness for windows across *Drosophila melanogaster* chromosome arms. In all plots, each point represents one window along the genome. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the two MDS coordinates against the midpoint of each window; rows correspond to chromosome arms. Colors are consistent for plots in each row. Vertical lines show the breakpoints of known polymorphic inversions. Solid black lines are for the inversions we used in Figure 3, while dotted grey lines are for other known inversions.

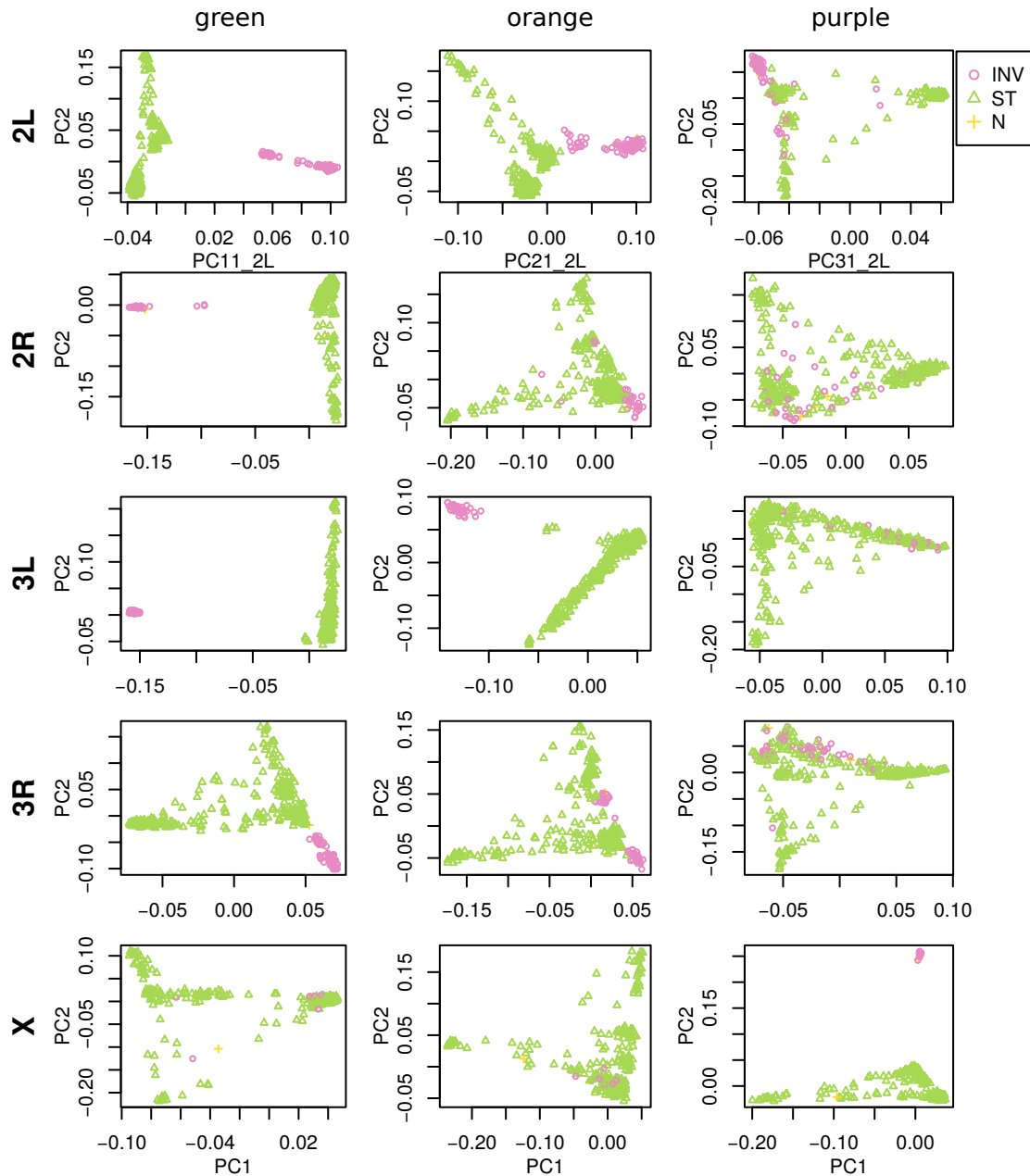


Figure 3: PCA plots for the three sets of genomic windows colored in Figure 2, on each chromosome arm of *Drosophila melanogaster*. In all plots, each point represents a sample. The first column shows the combined PCA plot for windows whose points are colored green in Figure 2; the second is for orange windows; and third is for purple windows. In each, samples are colored by orientation of the polymorphic inversions In(2L)t, In(2R)NS, In(3L)OK, In(3R)K and In(1)A respectively (data from [39]). In each “INV” denotes an inverted genotype, “ST” denotes the standard orientation, and “N” denotes unknown.

around the inversion breakpoints; there is also a relatively small dip on the right end of the arm (near the centromere). In contrast, Supplemental Figure S10 shows that after removing polymorphic inversions, remaining structure is dominated by the dip near the centromere. Without inversions, variation in patterns of relatedness shown in the MDS plots follows similar patterns to that previously seen in *D. melanogaster* recombination rate and diversity [40, 44]. Indeed, correlations between the recombination rate in each window and the position on the first MDS coordinate are highly significant (Spearman's  $\rho = 0.54$ ,  $p < 2 \times 10^{-16}$ ; Figures 4 and S11). This is consistent with the hypothesis that variation is due to selection, since the strength of linked selection increases with local gene density, measured in units of recombination distance. The number of genes – measured as the number of transcription start and end sites within each window – was not significantly correlated with MDS coordinate ( $p = 0.22$ ).

### 3.3 Human

As we did for the *Drosophila* data, we applied our method separately to all 22 human autosomes. On each, variation in patterns of relatedness was dominated by a small number of windows having similar patterns of relatedness to each other that differed dramatically from the rest of the chromosome. These may be primarily inversions: outlying windows coincide with three of the six large polymorphic inversions described in Antonacci et al. [1], notably a particularly large, polymorphic inversion on 8p23 (Figure 5). Similar plots for all chromosomes are shown in Supplementary Figures S12, S13, and S14. PCA plots of many outlying windows show a characteristic trimodal shape (shown for chromosome 8 in Figure S15), presumably distinguishing samples having each of the three diploid genotypes for each inversion orientation (although we do not have data on orientation status). This trimodal shape has been proposed as a method to identify inversions [43], but distinguishing this hypothesis from others, such as regions of low recombination rate, would require additional data.

We also applied the method on all 22 autosomes together, and found that, remarkably, the inversion on chromosome 8 is still the most striking outlying signal (Figure S16). Further investigation with a denser set of SNPs, allowing a finer genomic resolution, may yield other patterns.

### 3.4 *Medicago truncatula*

Unlike the other two species, the method applied separately on all eight chromosomes of *Medicago truncatula* showed similar patterns of gradual change in patterns of relatedness across each chromosome, with no indications of chromosome-specific patterns. This consistency suggests that the factor affecting the population structure for each chromosome is the same, as might be caused by varying strengths of linked selection. To verify that variation in the effects of population structure is shared across chromosomes, we applied

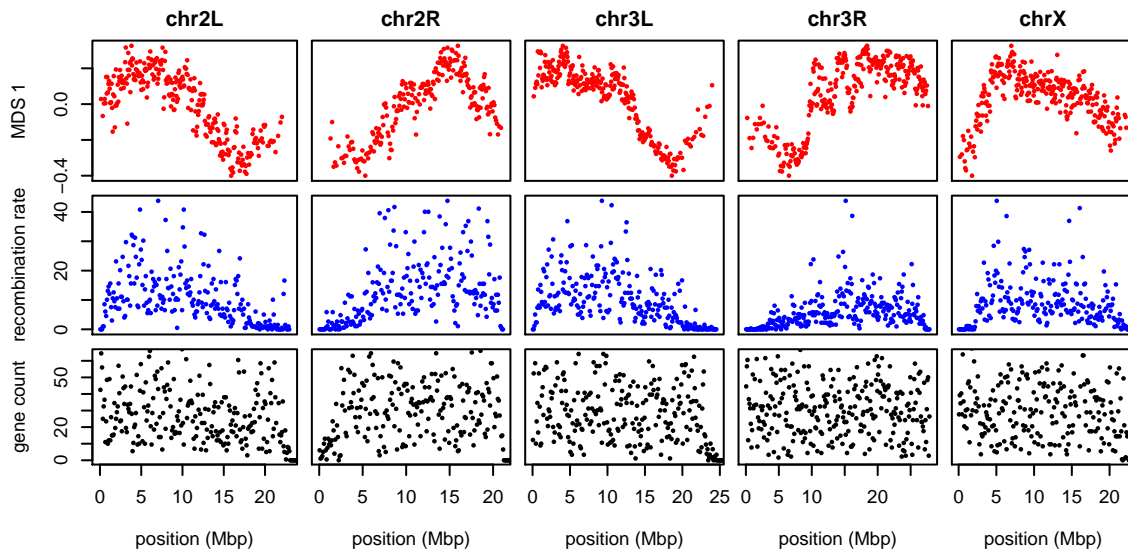


Figure 4: The effects of population structure without inversions is correlated to recombination rate in *Drosophila melanogaster*. The first plot (in red) shows the first MDS coordinate along the genome for windows of 10,000 SNPs, obtained after removing samples with inversions. (A plot analogous to Figure 2 is shown in Supplemental Figure S10.) The second plot (in blue) shows local average recombination rates in cM/Mbp, obtained as midpoint estimates for 100Kbp windows from the *Drosophila* recombination rate calculator [22] release 5, using rates from Comeron et al. [15]. The third plot (in black) shows the number of genes' transcription start and end sites within each 100Kbp window, divided by two. Transcription start and end sites were obtained from the RefGene table from the UCSC browser. The histone gene cluster on chromosome arm 2L is excluded.

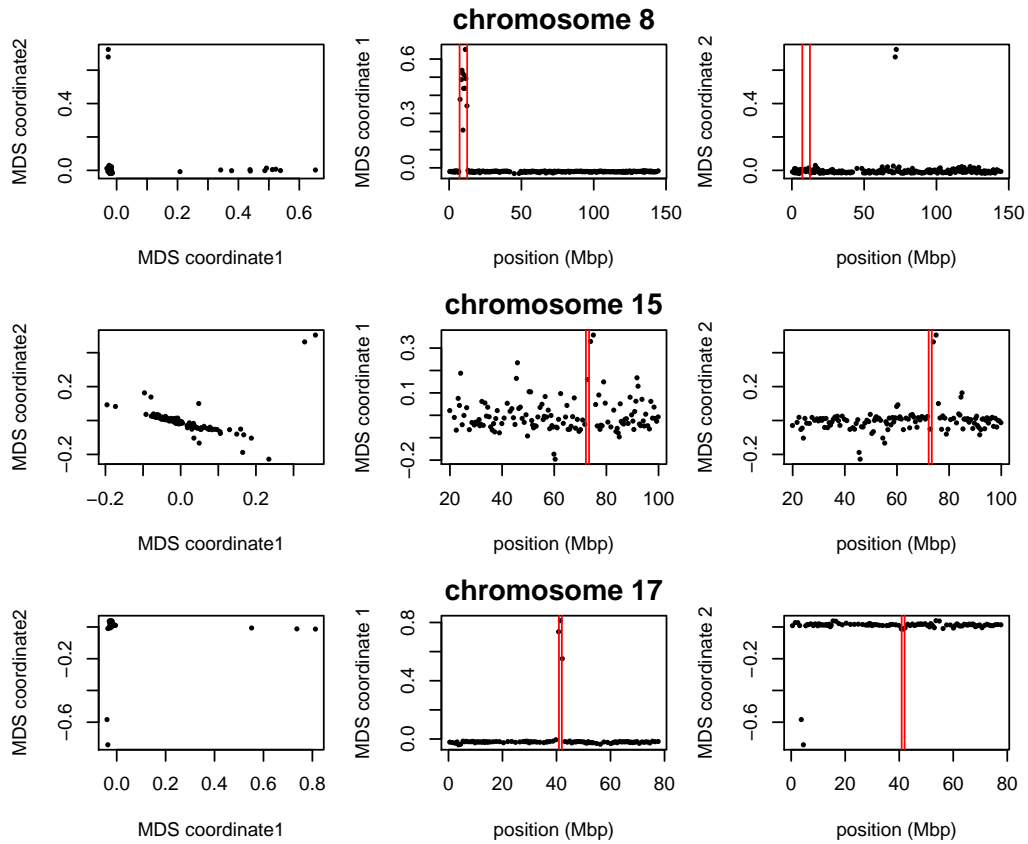


Figure 5: Variation in structure between windows on human chromosomes 8, 15, and 17. Each point in each plot represents a window. The first column shows the MDS visualization of relationships between windows; the second and third columns show the two MDS coordinates of each window against its position (midpoint) along the chromosome. Rows, from top to bottom show chromosomes 8, 15, and 17. The vertical red lines show the breakpoints of known inversions from Antonacci et al. [1].

the method to all chromosomes together. Results for chromosome 3 are shown in Figure 6, and other chromosomes are similar: across chromosomes, the high values of the first MDS coordinate coincide with the position of the heterochromatic regions surrounding the centromere, which often have lower gene density and may therefore be less subject to linked selection. To verify that this is a possible explanation, we counted the number of genes found in each window using gene models in Mt4.0 from [jcvi.org](http://jcvi.org) [63], which are shown juxtaposed with the first MDS coordinate of each window in Figure 7, and are significantly correlated, as shown in Supplemental Figure S17. (Values shown are the number of start and end positions of each predicted mRNA transcript, divided by two, assigned to the nearest window.) However, other genomic features, such as distance to centromere show roughly the same patterns, so we cannot rule out alternative hypotheses. In particular, fine-scale recombination rate estimates are not available in a form mappable to Mt4.0 coordinates (although those in Paape et al. [53] appear visually similar).

The results were highly consistent across window sizes, window types (SNPs or bp), and number of PCs, as shown in Table S2.

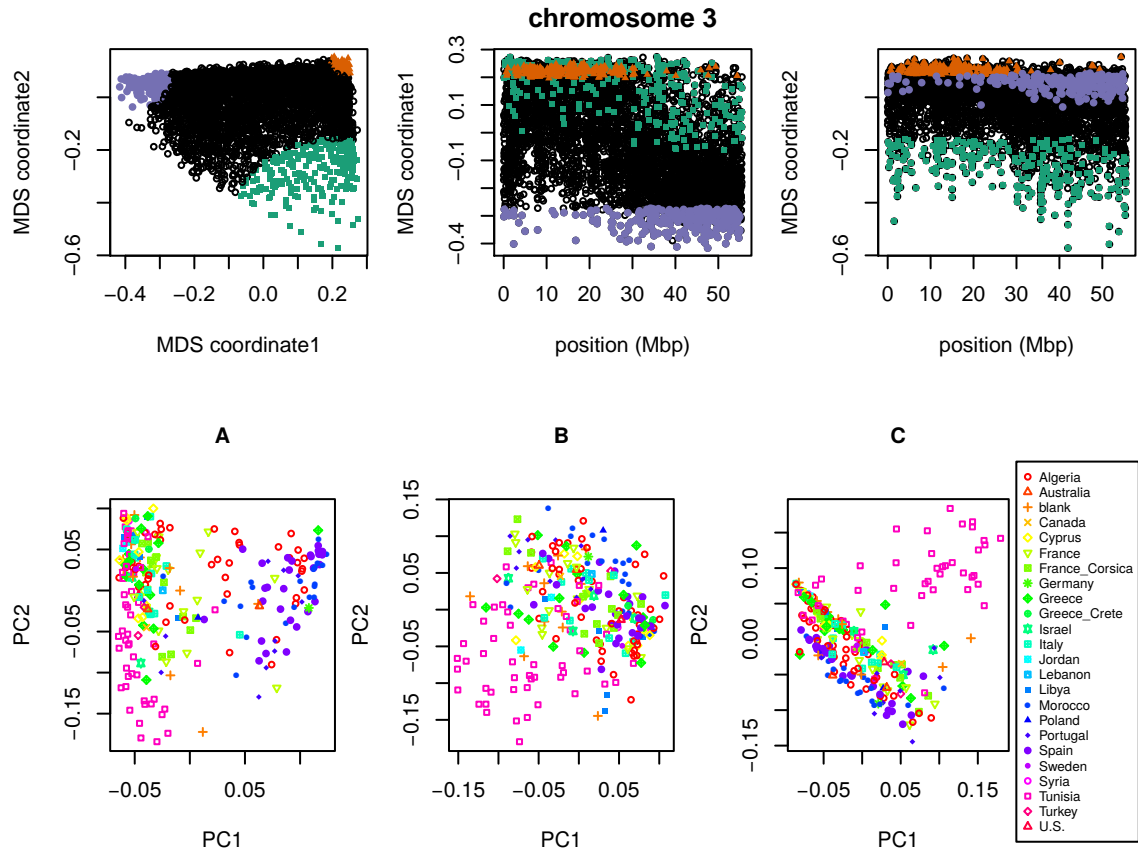
## 4 Discussion

Our investigations have found substantial variation in the patterns of relatedness formed by population structure across the genomes of three diverse species, revealing distinct biological processes driving this variation in each species. More investigation, particularly on more species and datasets, will help to uncover what aspects of species history can explain these differences. With growing appreciation of the heterogeneous effects of selection across the genome, especially the importance of adaptive introgression and hybrid speciation [5, 23, 31, 57, 62], local adaptation [41, 66], and inversion polymorphisms [37, 38], local PCA may prove to be a useful exploratory tool to discover important genomic features.

We now discuss possible implications of this variation in the effects of population structure, the impact of various parameter choices in implementing the method, and possible additional applications.

**Chromosomal inversions** A major driver of variation in patterns of relatedness in two datasets we examined seems to be inversions. This may be common, but the example of *Medicago truncatula* shows that polymorphic inversions are not ubiquitous. PCA has been proposed as a method for discovering inversions [43]; however, the signal left by inversions likely cannot be distinguished from long haplotypes under balancing selection or simply regions of reduced recombination without additional lines of evidence. Inversions show up in our method because across the inverted region, most gene trees share a common split that dates back to the origin of the inversion. However, in many applications, inversions are a nuisance. For instance, SMARTPCA [54] reduces their effect on PCA plots by regressing out the effect of linked SNPs on each other. Removing samples with the less common





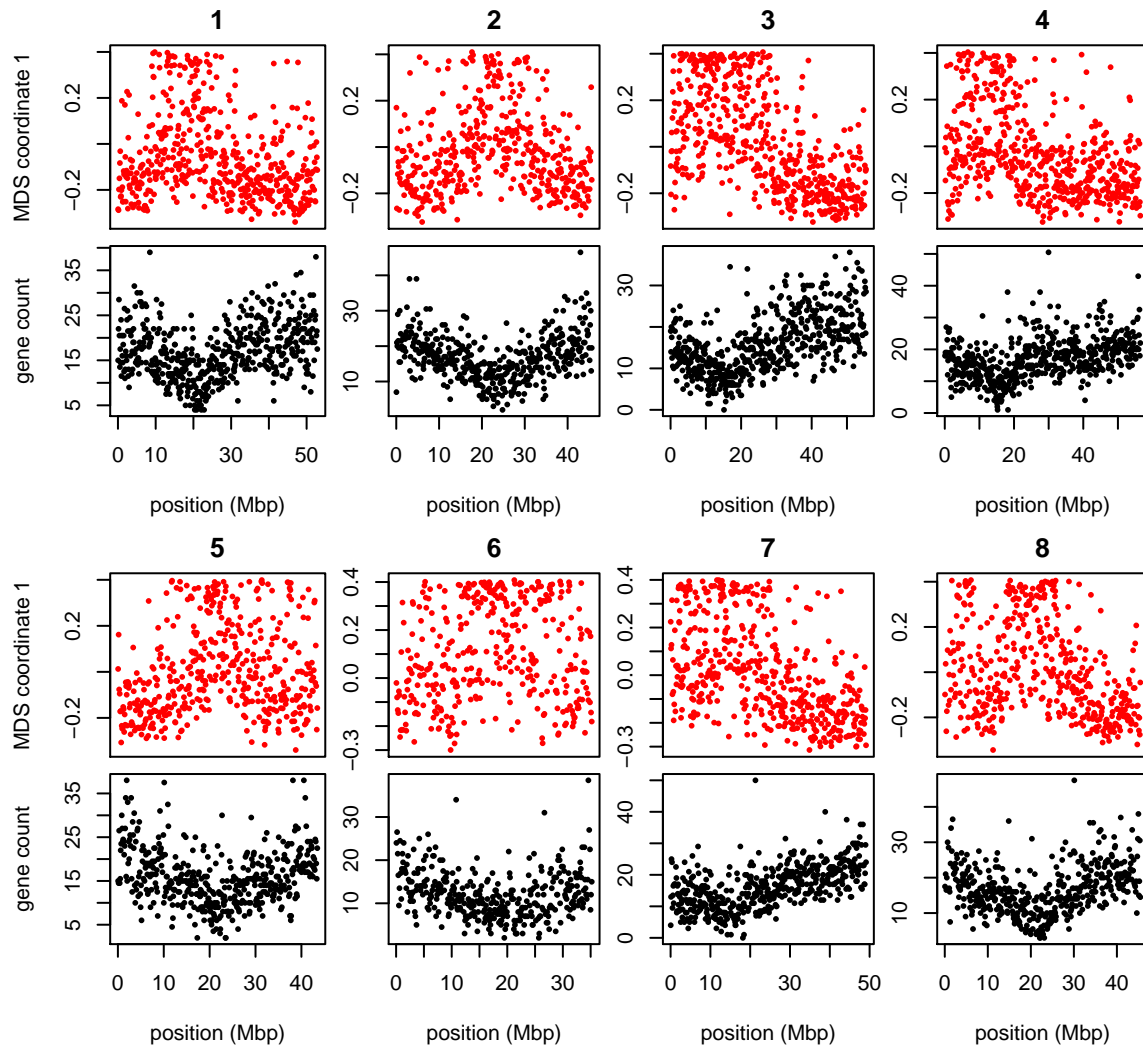


Figure 7: MDS coordinate and gene density for each window in the *Medicago* genome, for chromosomes 1–8 (numbered above each pair of figures). For each chromosome, the red plot above is first coordinate of MDS against the middle position of each window along each chromosome. The black plot below is gene count for each window against the middle position of each window.

orientation of each inversion reduced, but did not eliminate, the signal of inversions seen in the *Drosophila melanogaster* dataset, demonstrating that the genomic effects of transiently polymorphic inversions may outlast the inversions themselves.

**The effect of selection** Neutral processes are not expected to produce the chromosome-scale correlations we see in patterns of relatedness in the *Medicago truncatula* and *Drosophila melanogaster* datasets, because correlations induced by neutral processes should extend no further than does linkage disequilibrium (i.e., much less than a chromosome’s length). This suggests that they are produced by linked selection, a hypothesis backed up by correlations with gene density and recombination rate. We have also shown with simulations that linked selection can, in at least some circumstances, produce the sorts of patterns we observe. How might selection cause variation in patterns of relatedness? For instance, background selection (the effect on linked sites of selection against deleterious mutations Charlesworth et al. [10], Charlesworth [13]), can informally be thought of as reducing the number of potential contributors to the gene pool in regions of the genome with many possible deleterious mutations [29]. For this reason, if it acts in a spatial context, it is expected to induce samples from nearby locations to cluster together more frequently. Therefore, regions of the genome harboring many targets of local adaptation may show similar patterns, since migrant alleles in these regions will be selected against, and so locally gene trees will more closely reflect spatial proximity. Other forms of selection, such as hard sweeps on new mutations, repeated selection on standing variation, local adaptation, or temporally fluctuating selection, could clearly lead to variation in geographic patterns of relatedness in a similar way.

Another possible contributor is recent admixture between previously separated populations, the effects of which were not uniform across the genome due to selection. For instance, it has been hypothesized that large-scale variation in amount of introgressed Neanderthal DNA along the genome is due to selection against Neanderthal genes, leading to greater introgression in regions of lower gene density [27, 33]. African *Drosophila melanogaster* are thought to have a substantial amount of recently introgressed genome from “cosmopolitan” sources; if selection regularly favors genes from one origin, this could lead to substantial variation in patterns of relatedness correlated with local gene density.

There has been substantial debate over the relative impacts of different forms of selection (e.g., Burri et al. [8], Charlesworth et al. [11], Charlesworth [12], Corbett-Detig et al. [17], Harris and Nielsen [27], Hedrick [28], Martin et al. [46], Pease and Hahn [55], Phung et al. [56], Stankowski et al. [61]). These have been difficult to disentangle in part because for the most part theory makes predictions which are only strictly valid in randomly mating (i.e., unstructured) populations, and it is unclear to what extent the spatial structure observed in most real populations will affect these predictions. It may be possible to design more powerful statistics that make stronger use of spatial information.

**Parameter choices** There are several choices in the method that may in principle affect the results. As with whole-genome PCA, the choice of samples is important, as variation not strongly represented in the sample will not be discovered. The effects of strongly imbalanced sampling schemes are often corrected by dropping samples in overrepresented groups; but downweighting may be a better option that does not discard data. Next, the choice of window size may be important, although in our applications results were not sensitive to this. Finally, which collections of genomic regions are compared to each other (steps 3 and 4 in Figure 1), along with the method used to discover common structure, will affect results. We used MDS, applied to either each chromosome separately or to the entire genome; for instance, human inversions are clearly visible as outliers when compared to the rest of their chromosomes, but genome-wide, their signal is obscured by the numerous other signals of comparable strength.

Besides window length, there is also the question of how to choose windows. In these applications we have used nonoverlapping windows with equal numbers of polymorphic sites. However, we found little change in results when using different window sizes or when measuring windows in physical distance (in bp).

Finally, our software allows different choices for how many PCs to use in approximating structure of each window ( $k$  in equation 1), and how many MDS coordinates to use when describing the distance matrix between windows, but in our exploration, changing these has not produced dramatically different results. These are all part of more general techniques in dimension reduction and high-dimensional data visualization; we encourage the user to experiment.

**Applications** So-called cryptic relatedness between samples has been one of the major sources of confounding in genome-wide association studies (GWAS) and so methods must account for it by modeling population structure or kinship [2, 69]. Modern “mixed model” methods [e.g. 42] account for this with either a single, genome-wide kinship matrix or one constructed using only sites unlinked to the focal SNP. Since the effects of population structure is not constant along the genome, this could in principle lead to an inflation of false positives in parts of the genome with stronger population structure than the genome-wide average. A method such as ours might be used to estimate local kinship matrices, thus providing a more sensitive correction, although doing so without removing the signal itself could be challenging. Fortunately, in our human dataset this does not seem likely to have a strong effect: most variation is due to small, independent regions, possibly primarily inversions, and so may not have a major effect on GWAS. In the other species we examined, particularly *Drosophila melanogaster*, treating population structure as a single quantity would entail a substantial loss of power, and could potentially be misleading.

## Acknowledgements

We are indebted to John Pool, Russ Corbett-Detig, Matilde Cordeiro, and Peter Chang for assistance with obtaining data and interpreting results (especially inversion status of *D. melanogaster* samples). Jaime Ashander and Jerome Kelleher provided assistance in performing the simulations. Thanks also go to Yaniv Brandvain, Barbara Engelhardt, Charles Langley, Graham Coop, and Jeremy Berg for helpful comments and for encouraging the project.

## Disclosure declaration

The authors declare no conflicts of interest.

## References

- [1] Francesca Antonacci, Jeffrey M Kidd, Tomas Marques-Bonet, Mario Ventura, Priscillia Siswara, Zhaoshi Jiang, and Evan E Eichler. Characterization of six human disease-associated inversion polymorphisms. *Human molecular genetics*, 18(14):2555–2566, 2009.
- [2] William Astle and David J. Balding. Population structure and cryptic relatedness in genetic association studies. *Statistical Science*, 24(4):451–471, 11 2009. doi: 10.1214/09-STS307. URL <http://dx.doi.org/10.1214/09-STS307>.
- [3] Nicholas H. Barton. Genetic hitchhiking. *Philos Trans R Soc Lond B Biol Sci*, 355(1403):1553–1562, November 2000. doi: 10.1098/rstb.2000.0716. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1692896/>.
- [4] Albert P. Blair. Population structure in toads. *The American Naturalist*, 77(773):563–568, 1943. ISSN 00030147, 15375323. URL <http://www.jstor.org/stable/2457848>.
- [5] Yaniv Brandvain, Amanda M. Kenney, Lex Flagel, Graham Coop, and Andrea L. Sweigart. Speciation and introgression between *Mimulus nasutus* and *Mimulus guttatus*. *PLoS Genet*, 10(6):e1004410, 06 2014. doi: 10.1371/journal.pgen.1004410. URL <http://dx.doi.org/10.1371%2Fjournal.pgen.1004410>.
- [6] A. Brisbin, K. Bryc, J. Byrnes, F. Zakharia, L. Omberg, J. Degenhardt, A. Reynolds, H. Ostrer, J. G. Mezey, and C. D. Bustamante. PCAdmix: principal components-based assignment of ancestry along each chromosome in individuals with admixed ancestry from two or more populations. *Hum. Biol.*, 84(4):343–364, August 2012.
- [7] K Bryc, A Auton, M R Nelson, J R Oksenberg, S L Hauser, S Williams, A Froment, J M Bodo, C Wambebe, S A Tishkoff, and C D Bustamante. Genome-wide patterns

- of population structure and admixture in West Africans and African Americans. *Proc Natl Acad Sci U S A*, 107(2):786–791, January 2010. doi: 10.1073/pnas.0909559107. URL <http://www.ncbi.nlm.nih.gov/pubmed/20080753>.
- [8] R Burri, A Nater, T Kawakami, C F Mugal, P I Olason, L Smeds, A Suh, L Dutoit, S Bureš, L Z Garamszegi, S Hogner, J Moreno, A Qvarnström, M Ružić, S A Sæther, G P Sætre, J Török, and H Ellegren. Linked selection and recombination rate variation drive the evolution of the genomic landscape of differentiation across the speciation continuum of *Ficedula* flycatchers. *Genome Res*, 25(11):1656–1665, November 2015. doi: 10.1101/gr.196485.115. URL <https://www.ncbi.nlm.nih.gov/pubmed/26355005>.
- [9] Frank M.T.A. Busing, Erik Meijer, and Rien Van Der Leeden. Delete-m jackknife for unequal m. *Statistics and Computing*, 9(1):3–8, 1999. ISSN 0960-3174. doi: 10.1023/A:1008800423698. URL <http://dx.doi.org/10.1023/A%3A1008800423698>.
- [10] B Charlesworth, M T Morgan, and D Charlesworth. The effect of deleterious mutations on neutral molecular variation. *Genetics*, 134(4):1289–1303, August 1993. URL <http://www.genetics.org/content/134/4/1289>.
- [11] B Charlesworth, M Nordborg, and D Charlesworth. The effects of local selection, balanced polymorphism and background selection on equilibrium patterns of genetic diversity in subdivided populations. *Genet Res*, 70(2):155–174, October 1997. URL <https://www.ncbi.nlm.nih.gov/pubmed/9449192>.
- [12] Brian Charlesworth. The effects of deleterious mutations on evolution at linked sites. *Genetics*, 190(1):5–22, January 2012. doi: 10.1534/genetics.111.134288. URL <http://www.ncbi.nlm.nih.gov/pubmed/22219506>.
- [13] Brian Charlesworth. Background selection 20 years on: The Wilhelmine E. Key 2012 invitational lecture. *Journal of Heredity*, 104(2):161–171, 2013. doi: 10.1093/jhered/ess136. URL <http://jhered.oxfordjournals.org/content/104/2/161.abstract>.
- [14] Brian Charlesworth, Deborah Charlesworth, and Nicholas H. Barton. The effects of genetic and geographic structure on neutral variation. *Annual Review of Ecology, Evolution, and Systematics*, 34(1):99–125, 2003. doi: 10.1146/annurev.ecolsys.34.011802.132359. URL <http://arjournals.annualreviews.org/doi/abs/10.1146/annurev.ecolsys.34.011802.132359>.
- [15] J M Comeron, R Ratnappan, and S Bailin. The many landscapes of recombination in *Drosophila melanogaster*. *PLoS Genet*, 8(10), 2012. doi: 10.1371/journal.pgen.1002905. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3469467/>.
- [16] Russell B Corbett-Detig and Daniel L Hartl. Population genomics of inversion polymorphisms in *Drosophila melanogaster*. *PLoS Genet*, 8(12):e1003056, 2012.

- [17] Russell B. Corbett-Detig, Daniel L. Hartl, and Timothy B. Sackton. Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol*, 13(4):e1002112, 04 2015. doi: 10.1371/journal.pbio.1002112. URL <http://dx.doi.org/10.1371/journal.pbio.1002112>.
- [18] T E Cruickshank and M W Hahn. Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow. *Mol Ecol*, 23(13):3133–3157, 07 2014. doi: 10.1111/mec.12796. URL <https://www.ncbi.nlm.nih.gov/pubmed/24845075>.
- [19] Nicolas Duforet-Frebourg, Keurcien Luu, Guillaume Laval, Eric Bazin, and Michael G.B. Blum. Detecting genomic signatures of natural selection with principal component analysis: Application to the 1000 genomes data. *Molecular Biology and Evolution*, 2015. doi: 10.1093/molbev/msv334. URL <http://mbe.oxfordjournals.org/content/early/2016/01/12/molbev.msv334.abstract>.
- [20] B. Efron. *The Jackknife, the Bootstrap and Other Resampling Plans*. Society for Industrial and Applied Mathematics, 1982. doi: 10.1137/1.9781611970319. URL <http://epubs.siam.org/doi/abs/10.1137/1.9781611970319>.
- [21] Hans Ellegren, Linnea Smeds, Reto Burri, Pall I. Olason, Niclas Backstrom, Takeshi Kawakami, Axel Kunstner, Hannu Makinen, Krystyna Nadachowska-Brzyska, Anna Qvarnstrom, Severin Uebbing, and Jochen B. W. Wolf. The genomic landscape of species divergence in *Ficedula* flycatchers. *Nature*, 491(7426):756–760, November 2012. ISSN 00280836. doi: 10.1038/nature11584. URL <http://dx.doi.org/10.1038/nature11584>.
- [22] Anna-Sophie Fiston-Lavier, Nadia D. Singh, Mikhail Lipatov, and Dmitri A. Petrov. *Drosophila melanogaster* recombination rate calculator. *Gene*, 463(1–2):18 – 20, 2010. ISSN 0378-1119. doi: <http://dx.doi.org/10.1016/j.gene.2010.04.015>. URL <http://www.sciencedirect.com/science/article/pii/S0378111910001769>.
- [23] B M Fitzpatrick, J R Johnson, D K Kump, J J Smith, S R Voss, and H B Shaffer. Rapid spread of invasive genes into a threatened native species. *Proc Natl Acad Sci U S A*, 107(8):3606–3610, February 2010. doi: 10.1073/pnas.0911802107. URL <http://www.ncbi.nlm.nih.gov/pubmed/20133596>.
- [24] Rafael F. Guerrero, François Rousset, and Mark Kirkpatrick. Coalescent patterns for chromosomal inversions in divergent populations. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1587):430–438, 2011. ISSN 0962-8436. doi: 10.1098/rstb.2011.0246. URL <http://rstb.royalsocietypublishing.org/content/367/1587/430>.

- [25] Benjamin C. Haller and Philipp W. Messer. SLiM 2: Flexible, interactive forward genetic simulations. *Molecular Biology and Evolution*, 34(1):230–240, 2017. doi: 10.1093/molbev/msw211. URL [/brokenurl#+http://dx.doi.org/10.1093/molbev/msw211](http://dx.doi.org/10.1093/molbev/msw211).
- [26] Benjamin C. Haller, Jared Galloway, Jerome Kelleher, Philipp W. Messer, and Peter L. Ralph. Tree-sequence recording in SLiM opens new horizons for forward-time simulation of whole genomes. *bioRxiv*, 2018. doi: 10.1101/407783. URL <https://www.biorxiv.org/content/early/2018/09/04/407783>.
- [27] Kelley Harris and Rasmus Nielsen. The genetic cost of Neanderthal introgression. *Genetics*, 203(2):881–891, June 2016. URL <http://www.genetics.org/content/203/2/881>.
- [28] P W Hedrick. Adaptive introgression in animals: examples and comparison to new mutation and standing variation as sources of adaptive variation. *Mol Ecol*, 22(18):4606–4618, September 2013. doi: 10.1111/mec.12415. URL <https://www.ncbi.nlm.nih.gov/pubmed/23906376>.
- [29] R R Hudson and N L Kaplan. Deleterious background selection with recombination. *Genetics*, 141(4):1605–1617, December 1995. URL <http://www.genetics.org/content/141/4/1605>.
- [30] Emilia Huerta-Sánchez, Michael DeGiorgio, Luca Pagani, Ayele Tarekegn, Rosemary Ekong, Tiago Antao, Alexia Cardona, Hugh E. Montgomery, Gianpiero L. Cavalleri, Peter A. Robbins, Michael E. Weale, Neil Bradman, Endashaw Bekele, Toomas Kivisild, Chris Tyler-Smith, and Rasmus Nielsen. Genetic signatures reveal high-altitude adaptation in a set of Ethiopian populations. *Molecular Biology and Evolution*, 30(8):1877–1888, 2013. doi: 10.1093/molbev/mst089. URL <http://mbe.oxfordjournals.org/content/30/8/1877.abstract>.
- [31] Matthew B. Hufford, Pesach Lubinsky, Tanja Pyhäjärvi, Michael T. Devengenzo, Norman C. Ellstrand, and Jeffrey Ross-Ibarra. The genomic signature of crop-wild introgression in maize. *PLoS Genet*, 9(5):e1003477, 05 2013. doi: 10.1371/journal.pgen.1003477. URL <http://dx.doi.org/10.1371%2Fjournal.pgen.1003477>.
- [32] International HapMap Consortium, K A Frazer, D G Ballinger, D R Cox, D A Hinds, L L Stuve, R A Gibbs, J W Belmont, A Boudreau, P Hardenbol, S M Leal, S Pasternak, D A Wheeler, T D Willis, F Yu, H Yang, C Zeng, Y Gao, H Hu, W Hu, C Li, W Lin, S Liu, H Pan, X Tang, J Wang, W Wang, J Yu, B Zhang, Q Zhang, H Zhao, H Zhao, J Zhou, S B Gabriel, R Barry, B Blumenstiel, A Camargo, M Defelice, M Faggart, M Goyette, S Gupta, J Moore, H Nguyen, R C Onofrio, M Parkin, J Roy, E Stahl, E Winchester, L Ziaugra, D Altshuler, Y Shen, Z Yao, W Huang, X Chu, Y He, L Jin, Y Liu, Y Shen, W Sun, H Wang, Y Wang, Y Wang, X Xiong, L Xu,



- M M Waye, S K Tsui, H Xue, J T Wong, L M Galver, J B Fan, K Gunderson, S S Murray, A R Oliphant, M S Chee, A Montpetit, F Chagnon, V Ferretti, M Leboeuf, J F Olivier, M S Phillips, S Roumy, C Sallée, A Verner, T J Hudson, P Y Kwok, D Cai, D C Koboldt, R D Miller, L Pawlikowska, P Taillon-Miller, M Xiao, L C Tsui, W Mak, Y Q Song, P K Tam, Y Nakamura, T Kawaguchi, T Kitamoto, T Morizono, A Nagashima, Y Ohnishi, A Sekine, T Tanaka, T Tsunoda, P Deloukas, C P Bird, M Delgado, E T Dermitzakis, R Gwilliam, S Hunt, J Morrison, D Powell, B E Stranger, P Whittaker, D R Bentley, M J Daly, P I de Bakker, J Barrett, Y R Chretien, J Maller, S McCarroll, N Patterson, I Pe'er, A Price, S Purcell, D J Richter, P Sabeti, R Saxena, S F Schaffner, P C Sham, P Varilly, D Altshuler, L D Stein, L Krishnan, A V Smith, M K Tello-Ruiz, G A Thorisson, A Chakravarti, P E Chen, D J Cutler, C S Kashuk, S Lin, G R Abecasis, W Guan, Y Li, H M Munro, Z S Qin, D J Thomas, G McVean, A Auton, L Bottolo, N Cardin, S Eyheramendy, C Freeman, J Marchini, S Myers, C Spencer, M Stephens, P Donnelly, L R Cardon, G Clarke, D M Evans, A P Morris, B S Weir, T Tsunoda, J C Mullikin, S T Sherry, M Feolo, A Skol, H Zhang, C Zeng, H Zhao, I Matsuda, Y Fukushima, D R Macer, E Suda, C N Rotimi, C A Adebamowo, I Ajayi, T Aniagwu, P A Marshall, C Nkwodimmah, C D Royal, M F Leppert, M Dixon, A Peiffer, R Qiu, A Kent, K Kato, N Niikawa, I F Adewole, B M Knoppers, M W Foster, E W Clayton, J Watkin, R A Gibbs, J W Belmont, D Muzny, L Nazareth, E Sodergren, G M Weinstock, D A Wheeler, I Yakub, S B Gabriel, R C Onofrio, D J Richter, L Ziaugra, B W Birren, M J Daly, D Altshuler, R K Wilson, L L Fulton, J Rogers, J Burton, N P Carter, C M Clee, M Griffiths, M C Jones, K McLay, R W Plumb, M T Ross, S K Sims, D L Willey, Z Chen, H Han, L Kang, M Godbout, J C Wallenburg, P L'Archevêque, G Bellemare, K Saeki, H Wang, D An, H Fu, Q Li, Z Wang, R Wang, A L Holden, L D Brooks, J E McEwen, M S Guyer, V O Wang, J L Peterson, M Shi, J Spiegel, L M Sung, L F Zacharia, F S Collins, K Kennedy, R Jamieson, and J Stewart. A second generation human haplotype map of over 3.1 million SNPs. *Nature*, 449(7164):851–861, October 2007. doi: 10.1038/nature06258. URL <http://www.ncbi.nlm.nih.gov/pubmed/17943122>.
- [33] Ivan Juric, Simon Aeschbacher, and Graham Coop. The strength of selection against Neanderthal introgression. *bioRxiv*, 2016. doi: 10.1101/030148. URL <http://biorxiv.org/content/early/2016/07/22/030148>.
- [34] N. Kambhatla and T. K. Leen. Dimension reduction by local principal component analysis. *Neural Computation*, 9(7):1493–1516, July 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.7.1493. URL [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=6795533](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=6795533).
- [35] Jerome Kelleher, Kevin Thornton, Jaime Ashander, and Peter Ralph. Efficient pedigree recording for fast population genetics simulation. *bioRxiv*, 2018. doi: 10.1101/248500. URL <https://www.biorxiv.org/content/early/2018/06/07/248500>.

- [36] Yuseob Kim and Wolfgang Stephan. Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics*, 160(2):765–777, 2002. URL <http://www.genetics.org/cgi/content/abstract/160/2/765>.
- [37] Mark Kirkpatrick. How and why chromosome inversions evolve. *PLoS Biol*, 8(9), 2010. doi: 10.1371/journal.pbio.1000501. URL <http://www.ncbi.nlm.nih.gov/pubmed/20927412>.
- [38] Mark Kirkpatrick and Brian Barrett. Chromosome inversions, adaptive cassettes and the evolution of species’ ranges. *Molecular Ecology*, 2015. ISSN 1365-294X. doi: 10.1111/mec.13074. URL <http://dx.doi.org/10.1111/mec.13074>.
- [39] Justin B Lack, Charis M Cardeno, Marc W Crepeau, William Taylor, Russell B Corbett-Detig, Kristian A Stevens, Charles H Langley, and John E Pool. The *Drosophila* genome nexus: a population genomic resource of 623 *Drosophila melanogaster* genomes, including 197 from a single ancestral range population. *Genetics*, 199(4):1229–1241, 2015.
- [40] C H Langley, K Stevens, C Cardeno, Y C Lee, D R Schrider, J E Pool, S A Langley, C Suarez, R B Corbett-Detig, B Kolaczowski, S Fang, P M Nista, A K Holloway, A D Kern, C N Dewey, Y S Song, M W Hahn, and D J Begun. Genomic variation in natural populations of *Drosophila melanogaster*. *Genetics*, 192(2):533–598, October 2012. doi: 10.1534/genetics.112.142018. URL <http://www.ncbi.nlm.nih.gov/pubmed/22673804>.
- [41] Thomas Lenormand. Gene flow and the limits to natural selection. *Trends in Ecology & Evolution*, 17(4):183 – 189, 2002. ISSN 0169-5347. doi: DOI:10.1016/S0169-5347(02)02497-7. URL <http://www.sciencedirect.com/science/article/pii/S0169534702024977>.
- [42] P R Loh, G Tucker, B K Bulik-Sullivan, B J Vilhjálmsón, H K Finucane, R M Salem, D I Chasman, P M Ridker, B M Neale, B Berger, N Patterson, and A L Price. Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet*, 47(3):284–290, March 2015. doi: 10.1038/ng.3190. URL <https://www.ncbi.nlm.nih.gov/pubmed/25642633>.
- [43] J Ma and C I Amos. Investigation of inversion polymorphisms in the human genome using principal components analysis. *PLoS One*, 7(7), 2012. doi: 10.1371/journal.pone.0040224. URL <http://www.ncbi.nlm.nih.gov/pubmed/22808122>.
- [44] Trudy F. C. Mackay, Stephen Richards, Eric A. Stone, Antonio Barbadilla, Julien F. Ayroles, Dianhui Zhu, Sonia Casillas, Yi Han, Michael M. Magwire, Julie M. Cridland, Mark F. Richardson, Robert R. H. Anholt, Maite Barron, Crystal Bess, Kerstin Petra Blankenburg, Mary Anna Carbone, David Castellano, Lesley Chaboub, Laura Duncan,

- Zeke Harris, Mehwish Javid, Joy Christina Jayaseelan, Shalini N. Jhangiani, Katherine W. Jordan, Fremiet Lara, Faye Lawrence, Sandra L. Lee, Pablo Librado, Raquel S. Linheiro, Richard F. Lyman, Aaron J. Mackey, Mala Munidasa, Donna Marie Muzny, Lynne Nazareth, Irene Newsham, Lora Perales, Ling-Ling Pu, Carson Qu, Miquel Ramia, Jeffrey G. Reid, Stephanie M. Rollmann, Julio Rozas, Nehad Saada, Lavanya Turlapati, Kim C. Worley, Yuan-Qing Wu, Akihiko Yamamoto, Yiming Zhu, Casey M. Bergman, Kevin R. Thornton, David Mittelman, and Richard A. Gibbs. The *Drosophila melanogaster* genetic reference panel. *Nature*, 482(7384):173–178, February 2012. ISSN 00280836. URL <http://dx.doi.org/10.1038/nature10811>.
- [45] José V Manjón, Pierrick Coupé, Luis Concha, Antonio Buades, D Louis Collins, and Montserrat Robles. Diffusion weighted image denoising using overcomplete local PCA. *PloS one*, 8(9):e73021, 2013.
- [46] Simon Henry Martin, Markus Moest, Wiliam J Palmer, Camilo Salazar, W. Owen McMillan, Francis M Jiggins, and Chris D Jiggins. Natural selection and genetic diversity in the butterfly *Heliconius melpomene*. *Genetics*, 203(1):525–541, May 2016. doi: 10.1101/042796. URL <http://www.genetics.org/content/203/1/525>.
- [47] Gil McVean. A genealogical interpretation of principal components analysis. *PLoS Genet*, 5(10):e1000686, 2009.
- [48] P Menozzi, A Piazza, and L Cavalli-Sforza. Synthetic maps of human gene frequencies in Europeans. *Science*, 201(4358):786–792, September 1978. URL <http://www.ncbi.nlm.nih.gov/pubmed/356262>.
- [49] N J Nadeau, A Whibley, R T Jones, J W Davey, K K Dasmahapatra, S W Baxter, M A Quail, M Joron, R H French Constant, M L Blaxter, J Mallet, and C D Jiggins. Genomic islands of divergence in hybridizing *Heliconius* butterflies identified by large-scale targeted sequencing. *Philos Trans R Soc Lond B Biol Sci*, 367(1587):343–353, February 2012. doi: 10.1098/rstb.2011.0198. URL <http://www.ncbi.nlm.nih.gov/pubmed/22201164>.
- [50] M R Nelson, K Bryc, K S King, A Indap, A R Boyko, J Novembre, L P Briley, Y Maruyama, D M Waterworth, G Waeber, P Vollenweider, J R Oksenberg, S L Hauser, H A Stirnadel, J S Kooner, J C Chambers, B Jones, V Mooser, C D Bustamante, A D Roses, D K Burns, M G Ehm, and E H Lai. The Population Reference Sample, POPRES: a resource for population, disease, and pharmacological genetics research. *Am J Hum Genet*, 83(3):347–358, September 2008. doi: 10.1016/j.ajhg.2008.08.005. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2556436>.
- [51] John Novembre and Matthew Stephens. Interpreting principal component analyses of spatial population genetic variation. *Nature genetics*, 40(5):646–649, 2008.

- [52] John Novembre, Toby Johnson, Katarzyna Bryc, Zoltán Kutalik, Adam R Boyko, Adam Auton, Amit Indap, Karen S King, Sven Bergmann, Matthew R Nelson, et al. Genes mirror geography within Europe. *Nature*, 456(7218):98–101, 2008.
- [53] Timothy Paape, Peng Zhou, Antoine Branca, Roman Briskine, Nevin Young, and Peter Tiffin. Fine-scale population recombination rates, hotspots, and correlates of recombination in the *Medicago truncatula* genome. *Genome Biology and Evolution*, 4(5):726–737, 2012. doi: 10.1093/gbe/evs046. URL <http://gbe.oxfordjournals.org/content/4/5/726.abstract>.
- [54] Nick Patterson, Alkes L Price, and David Reich. Population structure and eigenanalysis. *PLoS Genetics*, 2(12):e190, 12 2006. doi: 10.1371/journal.pgen.0020190. URL <http://dx.plos.org/10.1371%2Fjournal.pgen.0020190>.
- [55] J B Pease and M W Hahn. More accurate phylogenies inferred from low-recombination regions in the presence of incomplete lineage sorting. *Evolution*, 67(8):2376–2384, August 2013. doi: 10.1111/evo.12118. URL <http://www.ncbi.nlm.nih.gov/pubmed/23888858>.
- [56] Tanya N. Phung, Christian D. Huber, and Kirk E. Lohmueller. Determining the effect of natural selection on linked neutral divergence across species. *PLOS Genetics*, 12(8):1–27, 08 2016. doi: 10.1371/journal.pgen.1006199. URL <https://doi.org/10.1371/journal.pgen.1006199>.
- [57] John E Pool. The mosaic ancestry of the *Drosophila* Genetic Reference Panel and the *D. melanogaster* reference genome reveals a network of epistatic fitness interactions. *Molecular Biology and Evolution*, 32(12):3236–3251, 2015. doi: 10.1101/014837. URL <http://mbe.oxfordjournals.org/content/32/12/3236.abstract>.
- [58] John E. Pool, Russell B. Corbett-Detig, Ryuichi P. Sugino, Kristian A. Stevens, Charis M. Cardeno, Marc W. Crepeau, Pablo Duchon, J. J. Emerson, Perot Saelao, David J. Begun, and Charles H. Langley. Population genomics of sub-Saharan *Drosophila melanogaster*: African diversity and non-African admixture. *PLoS Genet*, 8(12):1–24, 12 2012. doi: 10.1371/journal.pgen.1003080. URL <http://dx.doi.org/10.1371%2Fjournal.pgen.1003080>.
- [59] Alkes L Price, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick, and David Reich. Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics*, 38(8):904–909, 2006.
- [60] Sam T. Roweis and Lawrence K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. ISSN 0036-8075. doi: 10.1126/science.290.5500.2323. URL <http://science.sciencemag.org/content/290/5500/2323>.

- [61] Sean Stankowski, Madeline A Chase, Allison M Fuiten, Peter L Ralph, and Matthew A Streisfeld. The tempo of linked selection: Rapid emergence of a heterogeneous genomic landscape during a radiation of monkeyflowers. *bioRxiv*, 2018. doi: 10.1101/342352. URL <https://www.biorxiv.org/content/early/2018/06/21/342352>.
- [62] Fabian Staubach, Anna Lorenc, Philipp W. Messer, Kun Tang, Dmitri A. Petrov, and Diethard Tautz. Genome patterns of selection and introgression of haplotypes in natural populations of the house mouse (*Mus musculus*). *PLoS Genet*, 8(8):e1002891, 08 2012. doi: 10.1371/journal.pgen.1002891. URL <http://dx.doi.org/10.1371/journal.pgen.1002891>.
- [63] Haibao Tang, Vivek Krishnakumar, Shelby Bidwell, Benjamin Rosen, Agnes Chan, Shiguo Zhou, Laurent Genzbittel, Kevin L Childs, Mark Yandell, Heidrun Gundlach, et al. An improved genome release (version mt4. 0) for the model legume *Medicago truncatula*. *BMC genomics*, 15(1):1, 2014.
- [64] T L Turner, M W Hahn, and S V Nuzhdin. Genomic islands of speciation in *Anopheles gambiae*. *PLoS Biol*, 3(9), September 2005. doi: 10.1371/journal.pbio.0030285. URL <https://www.ncbi.nlm.nih.gov/pubmed/16076241>.
- [65] Benjamin Vernot and Joshua M. Akey. Resurrecting surviving neandertal lineages from modern human genomes. *Science*, 2014. doi: 10.1126/science.1245938. URL <http://www.sciencemag.org/content/early/2014/01/28/science.1245938.abstract>.
- [66] Ian J. Wang and Gideon S. Bradburd. Isolation by environment. *Molecular Ecology*, 23(23):5649–5662, 2014. ISSN 1365-294X. doi: 10.1111/mec.12938. URL <http://dx.doi.org/10.1111/mec.12938>.
- [67] Andreas Weingessel and Kurt Hornik. Local PCA algorithms. *Neural Networks, IEEE Transactions on*, 11(6):1242–1250, 2000.
- [68] Sewall Wright. The genetical structure of populations. *Annals of Eugenics*, 15(1): 323–354, 1949. ISSN 2050-1439. doi: 10.1111/j.1469-1809.1949.tb02451.x. URL <http://dx.doi.org/10.1111/j.1469-1809.1949.tb02451.x>.
- [69] Jian Yang, Noah A Zaitlen, Michael E Goddard, Peter M Visscher, and Alkes L Price. Advantages and pitfalls in the application of mixed-model association methods. *Nat Genet*, 46(2):100–106, February 2014. ISSN 10614036. URL <http://dx.doi.org/10.1038/ng.2876>.

## A Choosing window length

The choice of window length entails a balance between signal and noise. In very short windows, genealogies of the samples will only be represented by a few trees, so variation between windows represents demographic noise rather than meaningful variation in patterns of relatedness. Longer windows generally have more distinct trees (and SNPs), allowing for less noisy estimation of local patterns of relatedness. However, to better resolve meaningful signal, i.e., differences in patterns of relatedness along the genome, we would like reasonably short windows.

Since we summarize patterns of relatedness using relative positions in the principal component maps, we quantify “noise” as the standard error of a sample’s position on PC1 in a particular window, averaged across windows and samples, and “signal” as the standard deviation of the sample’s position on PC1 over all windows, averaged over samples. The definition of eigenvectors does not specify their sign, and so when comparing between windows we choose signs to best match each other: after choosing  $PC1_1$ , for instance, if  $u$  is the first eigenvector obtained from the covariance matrix for window  $j$ , then we next choose  $PC1_j = \pm u$ , where the sign is chosen according to which of  $\|PC1_1 - u\|$  or  $\|PC1_1 + u\|$  is smaller.

After doing this, the mean variance across windows is

$$\sigma_{\text{signal}}^2 = \frac{1}{N} \sum_{j=1}^N \frac{1}{L} \sum_{i=1}^L (PC1_{ij} - \overline{PC1}_j)^2,$$

where  $PC1_{ij}$  is the position of the  $i^{\text{th}}$  individual on PC1 in window  $j$ , and  $\overline{PC1}_j = (1/N) \sum_{i=1}^N PC1_{ij}$ . We estimate the standard error for each  $PC1_{ij}$  using the block jackknife [9, 20]: we divide the  $j^{\text{th}}$  window into 10 equal-sized pieces, and let  $PC1_{ij,k}$  denote the first principal component of this region found after removing the  $k^{\text{th}}$  piece; then the estimate of the squared standard error is  $\sigma_{ij}^2 = \frac{9}{10} \sum_{k=1}^{10} (PC1_{ij,k} - \frac{1}{10} \sum_{\ell=1}^{10} PC1_{ij,\ell})^2$ . Averaging over samples and windows,

$$\sigma_{\text{noise}}^2 = \frac{1}{N} \sum_{j=1}^N \frac{1}{L} \sum_{i=1}^L \sigma_{ij}^2.$$

For the main analysis, we defined windows to each consist of the same number of neighboring SNPs, and calculated  $\sigma_{\text{signal}}^2$  and  $\sigma_{\text{noise}}^2$  for a range of window sizes (i.e., numbers of SNPs). For our main results we chose the smallest window for which  $\sigma_{\text{signal}}^2$  was consistently larger than  $\sigma_{\text{noise}}^2$  (but checked other sizes); the values for various window sizes across *Drosophila* chromosomes are shown in Table S1. In the cases we examined, we found nearly identical results after varying window size, and choosing windows to be of the same physical length (in bp) rather than in numbers of SNPs.

chrom.	arm		window length (SNPs)				
			100	500	1,000	10,000	100,000
2L		$\sigma_{\text{noise}}^2$	2.05	1.64	1.18	0.17	0.04
		$\sigma_{\text{signal}}^2$	2.76	2.69	2.23	0.68	0.31
2R		$\sigma_{\text{noise}}^2$	2.18	1.92	1.63	0.58	0.13
		$\sigma_{\text{signal}}^2$	2.78	2.70	2.65	2.31	1.82
3L		$\sigma_{\text{noise}}^2$	2.08	2.00	1.64	0.73	0.25
		$\sigma_{\text{signal}}^2$	2.60	2.52	2.40	1.68	1.89
3R		$\sigma_{\text{noise}}^2$	1.95	1.76	1.44	0.59	0.20
		$\sigma_{\text{signal}}^2$	2.58	2.51	2.44	1.96	1.40
X		$\sigma_{\text{noise}}^2$	2.48	2.04	1.54	1.62	0.17
		$\sigma_{\text{signal}}^2$	2.61	2.43	2.30	0.32	1.14

Table S1: Measures of signal and noise, computed separately for each chromosome arm in the *Drosophila* dataset, at different window sizes. All values are multiplied by 1,000. Starting at windows of 1,000 SNPs, the signal (variation of PC1 between windows) starts to be substantially larger than the noise (standard error of PC1 for each window).

## B Simulations

We implemented two types of simulation: first, simple simulations of Gaussian “genotypes” where the expectation of variation in “population structure” was clear; and next, individual-based simulations with explicit genomes, using SLiM.

### B.1 Gaussian simulations

We simulated genotypes at each locus independently, drawing each vector of genotypes from a multivariate Gaussian distribution with zero mean and covariance matrix  $\Sigma$ . Sampled individuals came from three populations, and each  $\Sigma_{ij}$  depends on which populations the individuals  $i$  and  $j$  are in, as well as the location along the chromosome. There are three population-level mean relatedness matrices along the genome, which apply to the first

quarter ( $S^{(1)}$ ), the middle half ( $S^{(2)}$ ), and the last quarter ( $S^{(3)}$ ), respectively:

$$S^{(1)} = \begin{bmatrix} 0.75 & 0.25 & 0.0 \\ 0.25 & 0.75 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix}$$
$$S^{(2)} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 0.75 & 0.25 \\ 0.0 & 0.25 & 0.75 \end{bmatrix}$$
$$S^{(3)} = \begin{bmatrix} 0.75 & 0.0 & 0.25 \\ 0.0 & 1.0 & 0.0 \\ 0.25 & 0.0 & 0.75 \end{bmatrix}$$

If individuals  $i$  and  $j$  are in populations  $p(i)$  and  $p(j)$  respectively, then the covariance between their genotypes is  $\Sigma_{ij} = S_{p(i),p(j)}$ , using the appropriate  $S$  for that segment of the genome. The variance of individual  $i$ 's genotype is  $\Sigma_{ii} = S_{p(i),p(i)} + 0.1$ .

We first created “genotypes” in this way with fifty individuals from each of the three populations; running our method on a genome with 99 windows of 400 loci each produced the first plot in Figure S4. These matrices are chosen so that the top two eigenvalues  $\Sigma$  are the same (both 50.1), and so the ordering of the top two PCs is arbitrary. If our method was sensitive to PC ordering, then half the windows in each region that have one ordering would cluster with each other, separate from the other half.

We then marked each genotype in the first half of the chromosome as missing, independently, with probability 1/2 and ran our method again, producing the second plot of Figure S4. If our method was influenced by missing data, we would expect the first half of the chromosome to separate from the second in the MDS plot.

## B.2 SLiM simulations

Our SLiM simulations were constructed as follows. Individuals are diploid, and genomes have a length of 153,520,244 bp. Recombination was either (a) flat, with a constant rate of  $10^{-9}$ ; (b) according to the human female HapMap map for chromosome 7; or (c) constant in each of seven equal-sized regions, beginning at  $2.04 \times 10^{-8}$ , descending by a factor of four for three steps, and then ascending by a factor of four for three steps, so that the middle seventh has the lowest recombination rate, and the outer two sevenths has a rate 64 times higher. Selected mutations are introduced at a rate of  $10^{-10}$  per bp per individual per generation, and have selection coefficients drawn from a Gamma distribution with mean 0.005 and shape 2; each coefficient are either positive or negative with probabilities 1/30 and 29/30 respectively. Each simulation was run for 50,000 generations.

Each individual has a spatial position in the two-dimensional square of width  $W = 8$ . Each time step, each individual chooses the nearest other to mate with, producing a



random, Poisson distributed number of offspring with mean  $1/3$ . Offspring are assigned random spatial locations displaced from their parent's by a bivariate Gaussian with mean zero and standard deviation  $\sigma = 0.2$ , reflected to stay within the habitat range.

Each individual survives to the next time step with probability equal to their fitness. Fitness values are determined multiplicatively by the effects of each mutation, but are multiplied by an additional factor determined by the local density of individuals. This factor is equal to  $\rho/(1 + C)$ , where  $\rho = 2\pi K\sigma^2$  is the carrying capacity per circle of radius  $\sigma$ ;  $K = 100$  is the mean equilibrium population density; and  $C$  is the sum of a Gaussian kernel with standard deviation  $\sigma = 0.1$  between the focal individual and all other individuals within distance  $3\sigma$ . To avoid edge effects, fitnesses are further multiplied by  $\min(1, z)$ , where  $z$  is the distance to the nearest boundary. This produces populations that fluctuate at equilibrium around 6,000 individuals in total, fairly evenly spread across the square.

In one additional simulation, we modified fitnesses by multiplying the selective effect of each allele in each individual by multiplying it by  $2x/W - 1$ , where  $x$  is the  $x$  coordinate of the individual. This makes the effect of each allele opposite on the left and on the right, and neutral in the middle, and leads to a moderate number of balanced polymorphisms.

## C Supplementary Tables

## D Supplementary Figures

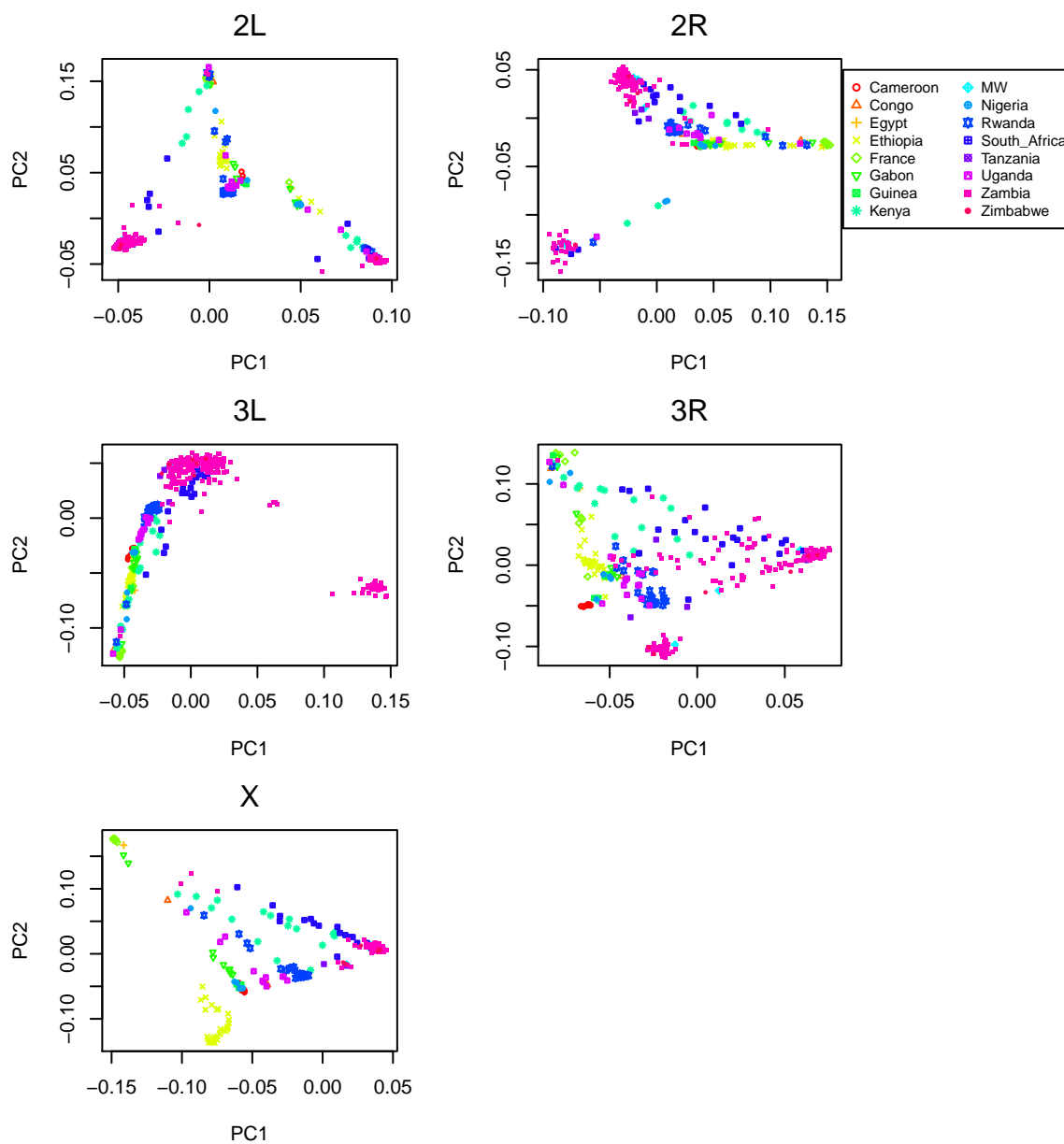


Figure S1: PCA plots for chromosome arms 2L, 2R, 3L, 3R and X of the *Drosophila melanogaster* dataset.

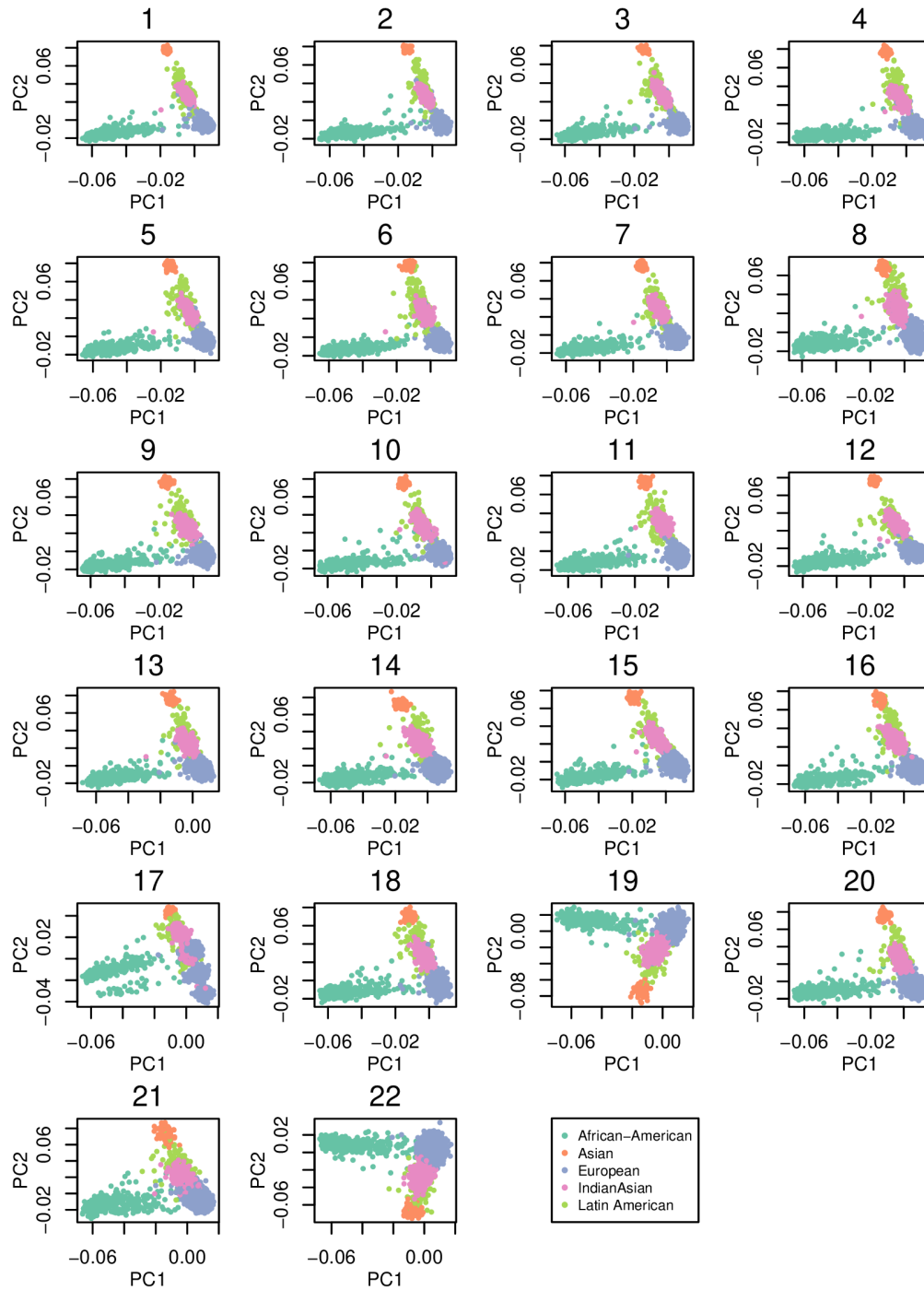


Figure S2: PCA plots for all 22 human autosomes from the POPRES data.

	10000 SNPs	1000 SNPs	10000 SNPs	100000bp	10000bp
MDS1	2 PCs	2 PCs	5 PCs	2 PCs	2 PCs
10000 SNPs, 2 PCs	1.00	0.87	0.96	0.90	0.88
1000 SNPs, 2 PCs	0.68	1.00	0.73	0.68	0.94
10000 SNPs, 5 PCs	0.96	0.92	1.00	0.88	0.93
100000bp, 2 PCs	0.90	0.87	0.88	1.00	0.87
10000bp, 2 PCs	0.68	0.93	0.72	0.67	1.00
MDS2					
10000 SNPs, 2 PCs	1.00	0.54	0.93	0.87	0.56
1000 SNPs, 2 PCs	0.82	1.00	0.76	0.83	0.92
10000 SNPs, 5 PCs	0.93	0.50	1.00	0.83	0.52
100000bp, 2 PCs	0.87	0.59	0.84	1.00	0.58
10000bp, 2 PCs	0.83	0.92	0.77	0.84	1.00

Table S2: Correlations between MDS coordinates of genomic regions between runs with different parameter values. To produce these, we first ran the algorithm with the specified window size and number of PCs ( $k$  in equation (1)) on the full *Medicago truncatula* dataset. Then to obtain the correlation between results obtained from parameters A in the row of the matrix above and parameters B in the column of the matrix above, we mapped the windows of B to those of A by averaging MDS coordinates of any windows of B whose midpoints lay in the corresponding window of A; we then computed the correlation between the MDS coordinates of A and the averaged MDS coordinates of B. This is not a symmetric operation, so these matrices are not symmetric. As expected, parameter values with smaller windows produce noisier estimates, but plots of MDS values along the genome are visually very similar.

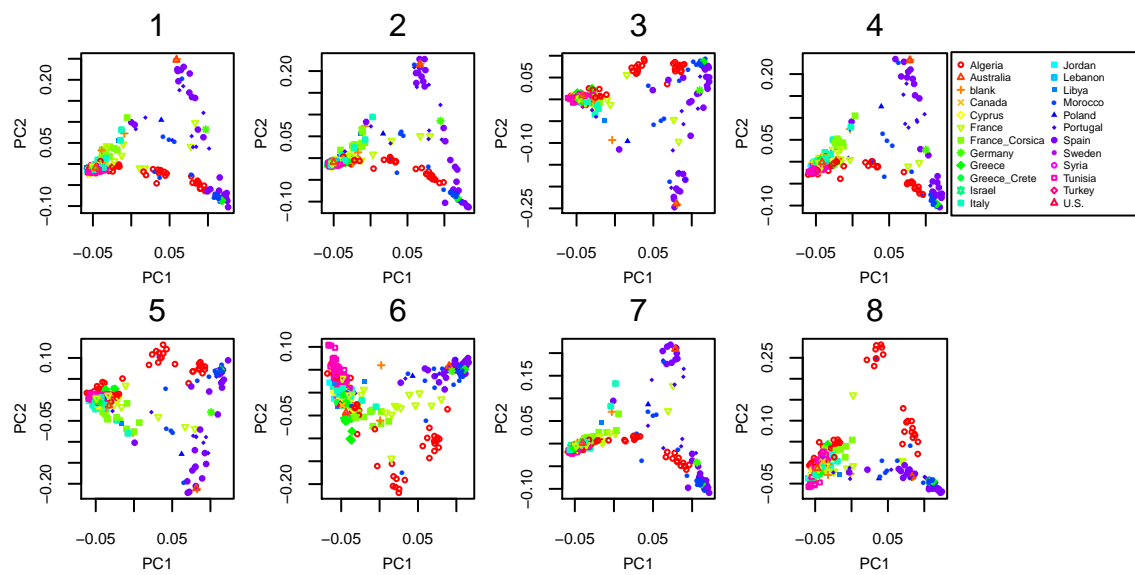


Figure S3: PCA plots for all 8 chromosomes in the *Medicago truncatula* dataset.

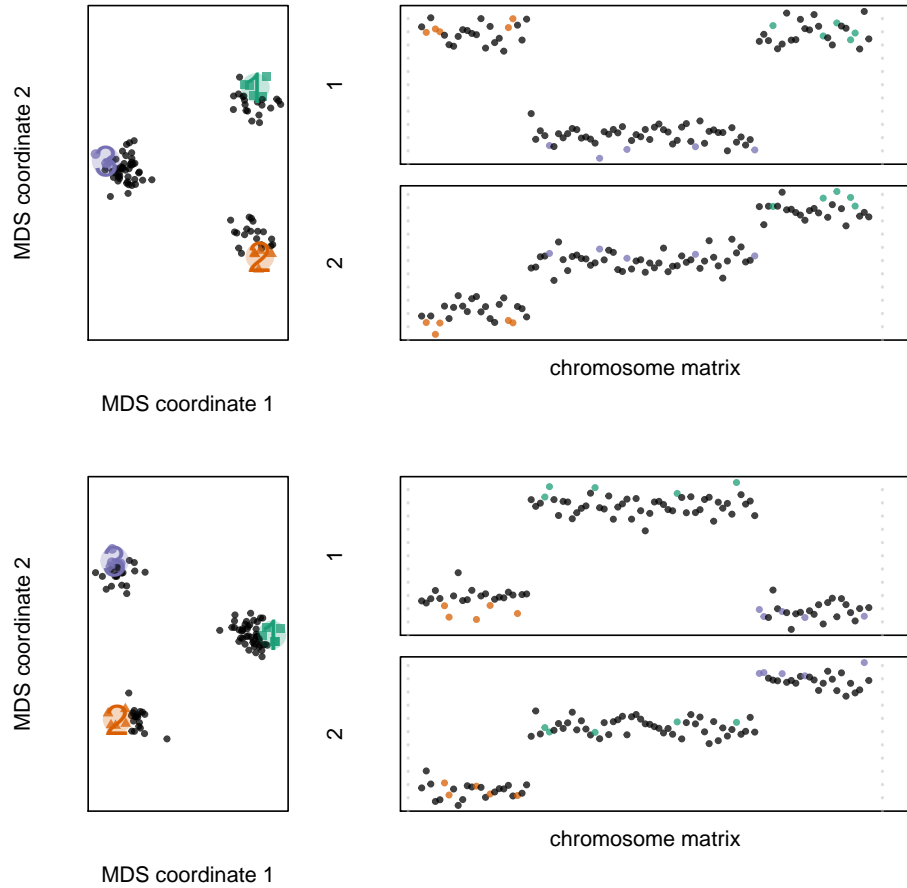


Figure S4: MDS visualizations of the Gaussian genotypes described in Appendix B.1, for 50 individuals from each of three populations. **(top)** The first quarter, middle half, and final quarter of the chromosome each have different population structure, as expected, despite the possibility for PC switching within each. **(bottom)** The same picture results even after marking a random 50% of the genotypes in the first half of the chromosome as missing.

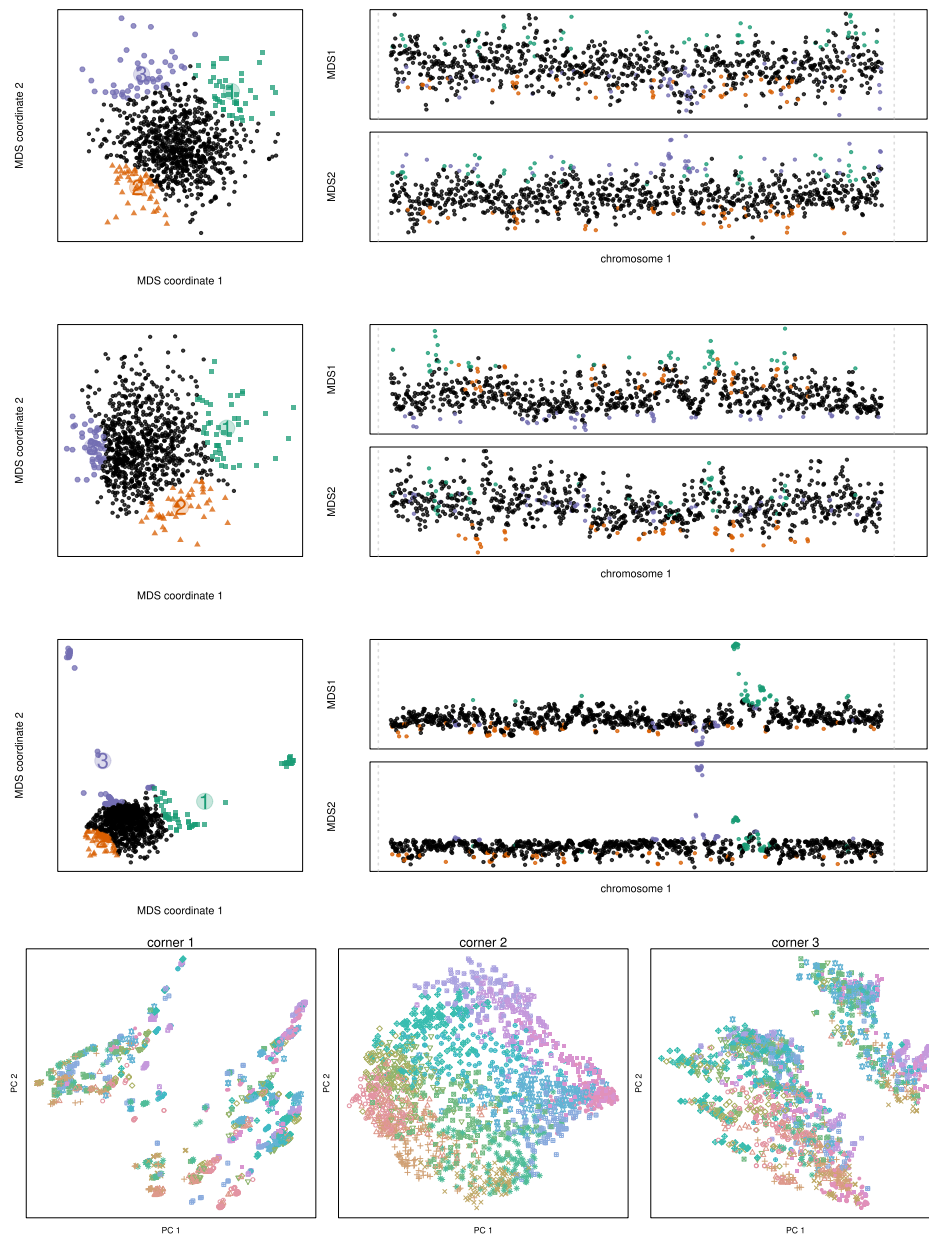


Figure S5: MDS visualizations of the results of individual-based simulations using SLiM (see Appendix B.2 for details). All simulations are neutral, and recombination is: **(top)** constant; **(top middle)** varies stepwise by factors of two in seven equal-length segments, with highest rates on the ends, so the middle segment has a recombination rate 64 times lower than the ends; **(bottom middle)** according to the HapMap human female chromosome 7 map. The **bottom** figure shows PCA maps corresponding to the three colored windows of the last (HapMap) situation; the<sup>39</sup> outlying regions are long regions of low recombination rate, so that region can be dominated by a few correlated trees, similar to an inversion.

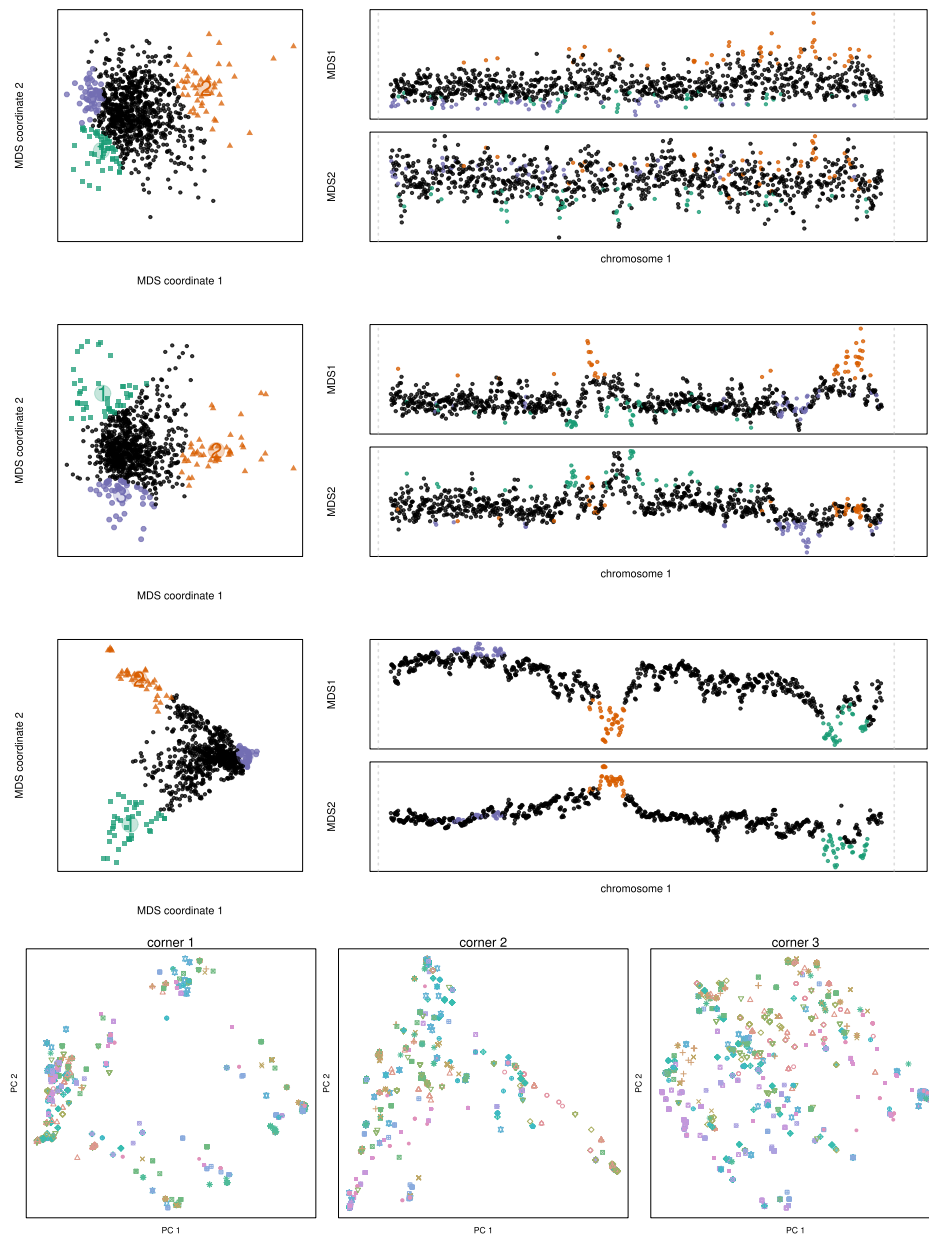


Figure S6: MDS visualizations of the results of individual-based simulations using SLiM (see Appendix B.2 for details). All simulations incorporate linked selection by allowing selected mutations to appear in the same two regions of the genome: the one-sixth of the genome immediately before the halfway point, and the last one-sixth of the genome. **(top)** Constant recombination rate. **(top middle)** Stepwise varying recombination rate (as described in Figure S5). **(bottom middle)** Constant recombination rate with spatially varying effects of selection. **(bottom)** PCA plots corresponding to the highlighted corners of the last MDS visualization, showing how spatially varying linked selection has affected patterns of relatedness.



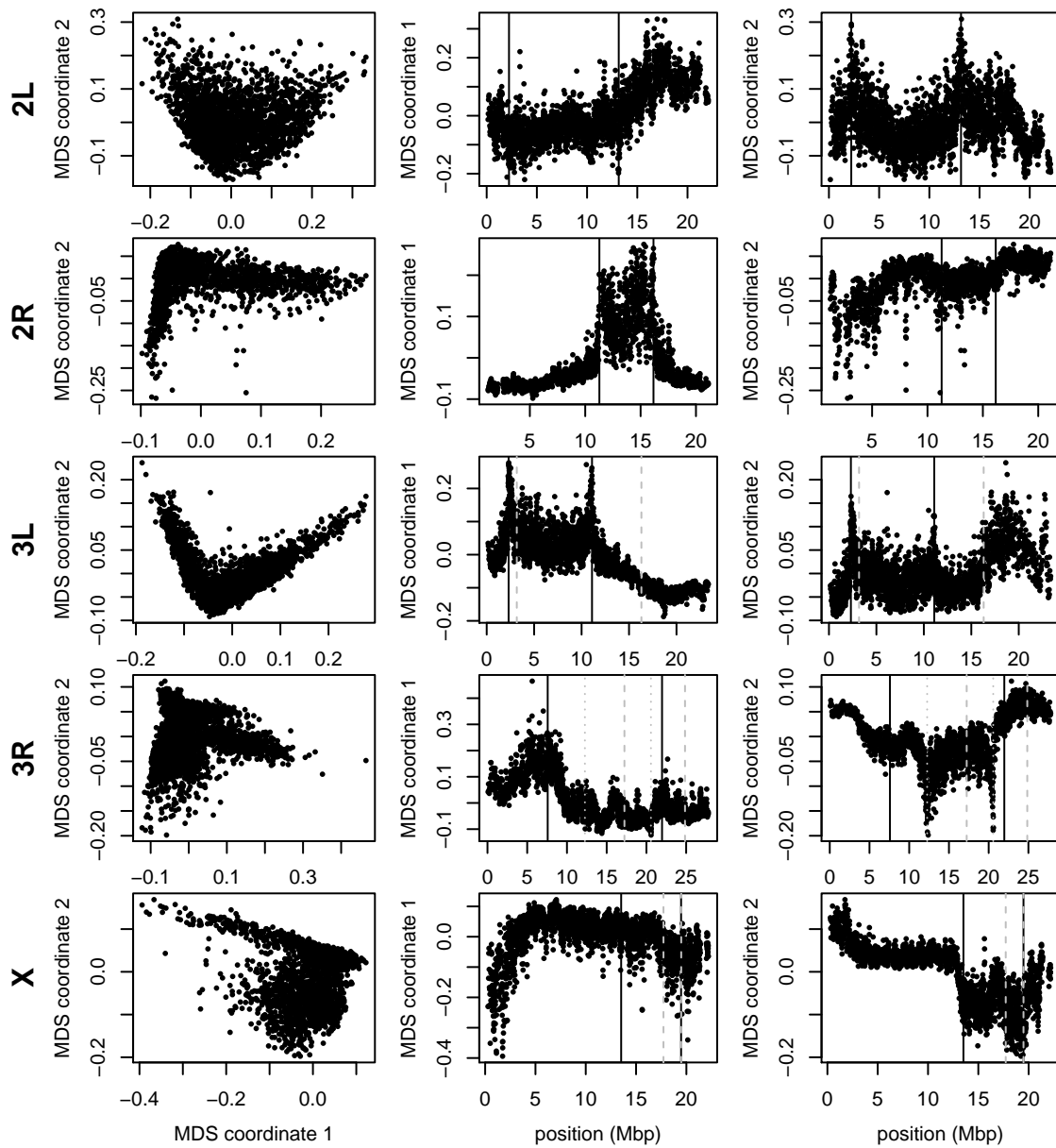


Figure S7: MDS visualizations for each chromosome arm of *Drosophila melanogaster*, as in Figure 2, except that the method was run using five PCs ( $k = 5$ ) instead of two.

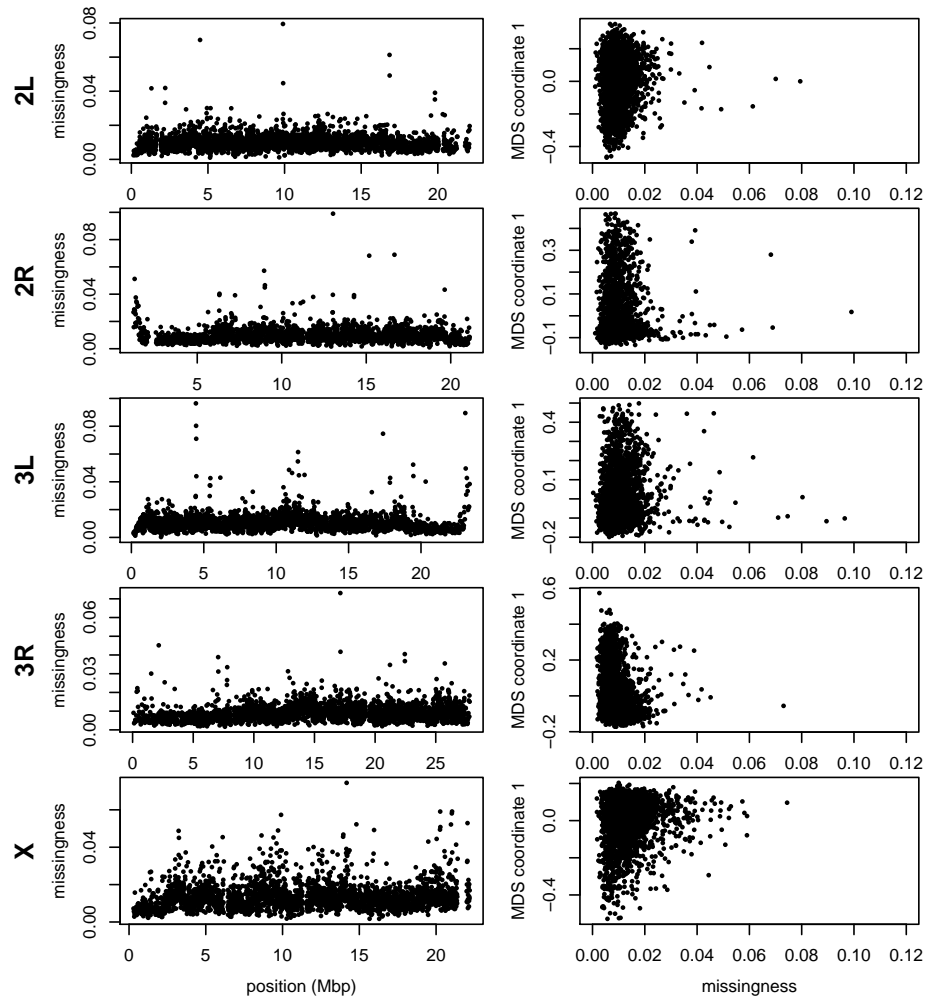


Figure S8: The proportion of data in each window that are missing, compared to the value of the first MDS coordinate for the *Drosophila melanogaster* data from Figure 2.

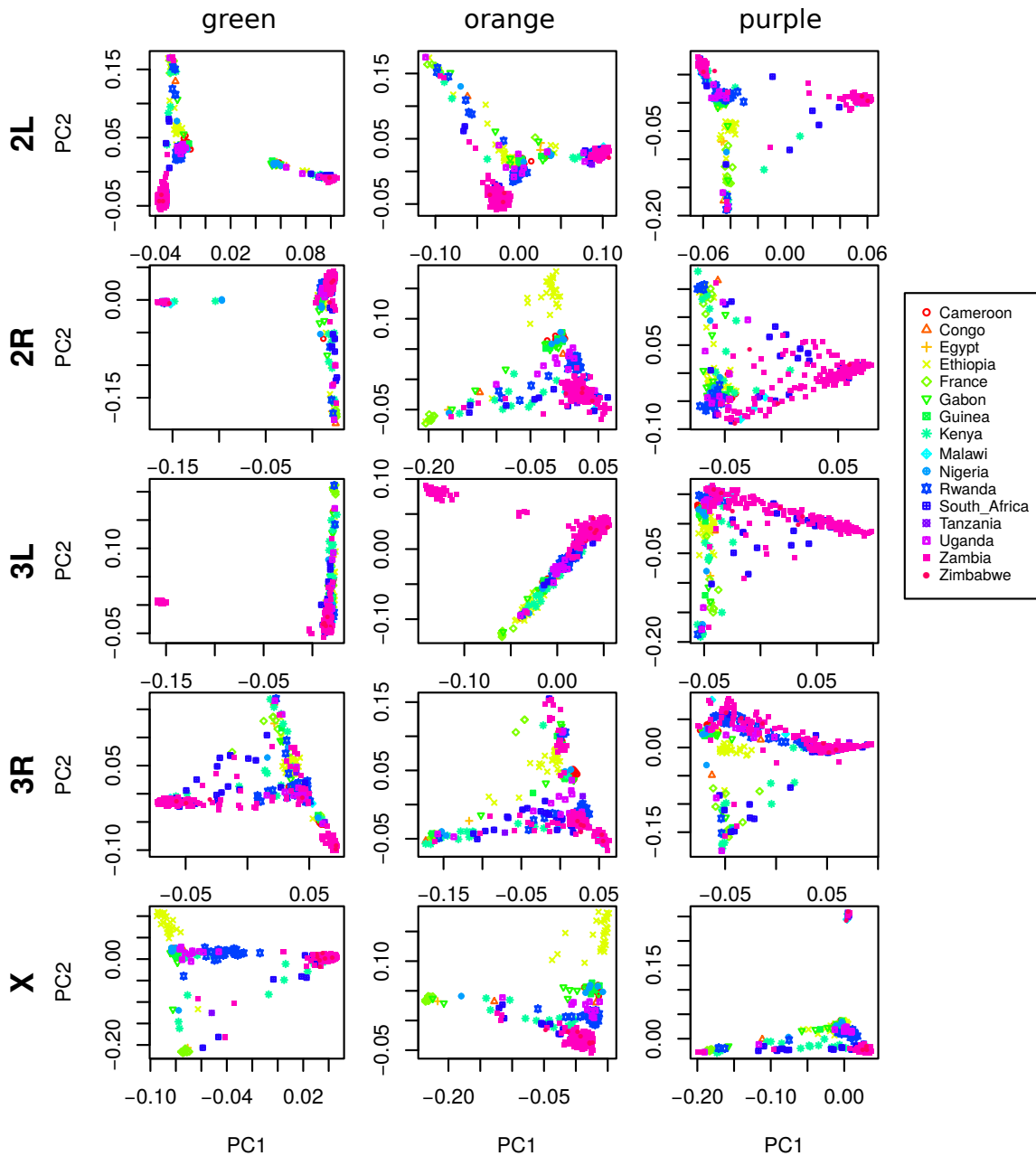


Figure S9: PCA plots for the three sets of genomic windows colored in Figure 2, on each chromosome arm of *Drosophila melanogaster*. In all plots, each point represents a sample. The first column shows the combined PCA plot for windows whose points are colored green in Figure 2; the second is for orange windows; and the third is for purple windows.

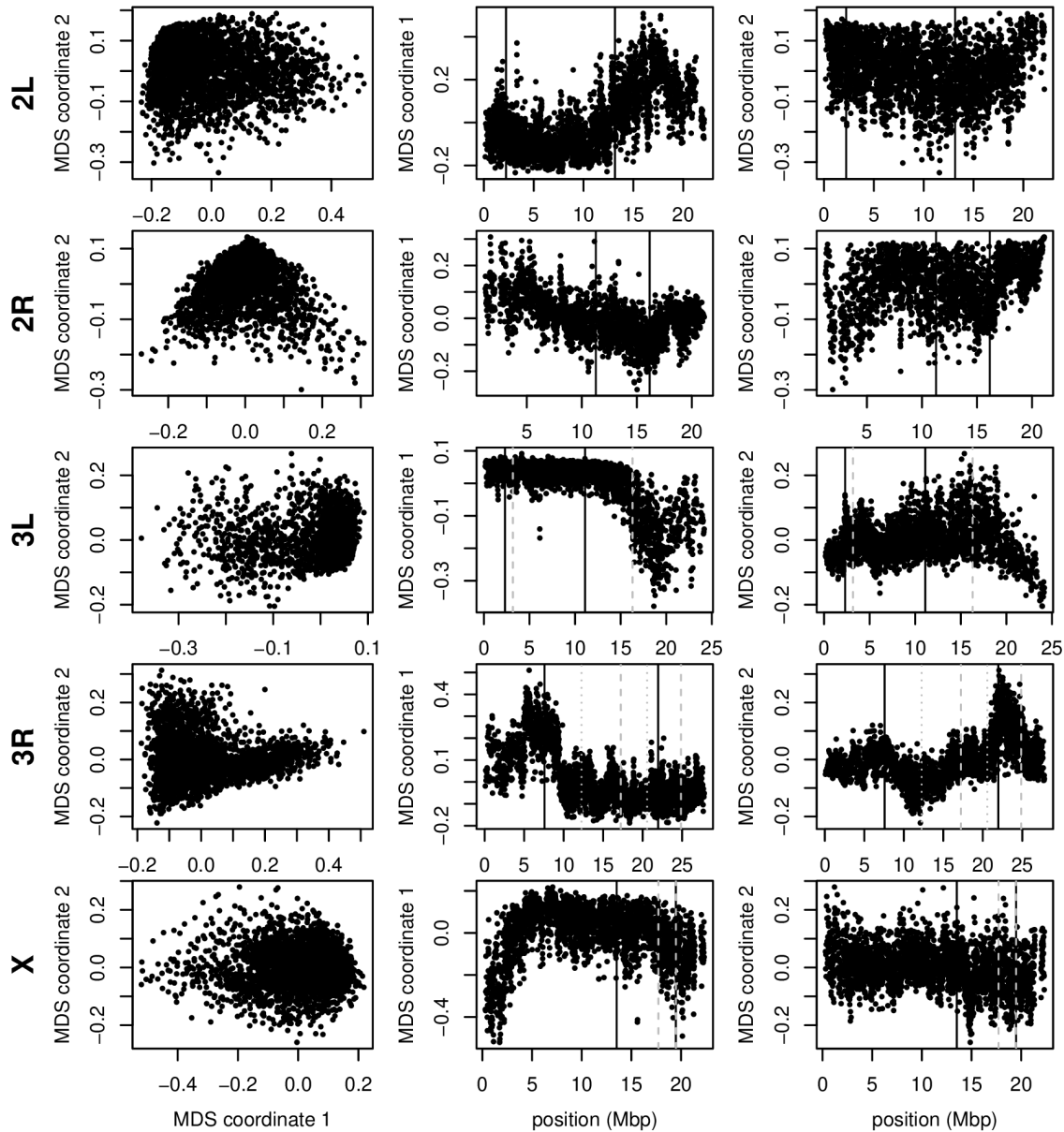


Figure S10: Variation in structure for windows of 1,000 SNPs across *Drosophila melanogaster* chromosome arms: without inversions. As in Figure 2, but after omitting for each chromosome arm individuals carrying the less frequent orientation of any inversions on that chromosome arm. The values differ from those in 4 in the window size used and that some MDS values were inverted (but relative orientation is meaningless as chromosome arms were run separately, unlike for *Medicago*). In all plots, each point represents one window along the genome. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the midpoint of each window against the two MDS coordinates; rows correspond to chromosome arms. Vertical lines show the breakpoints of known polymorphic inversions.

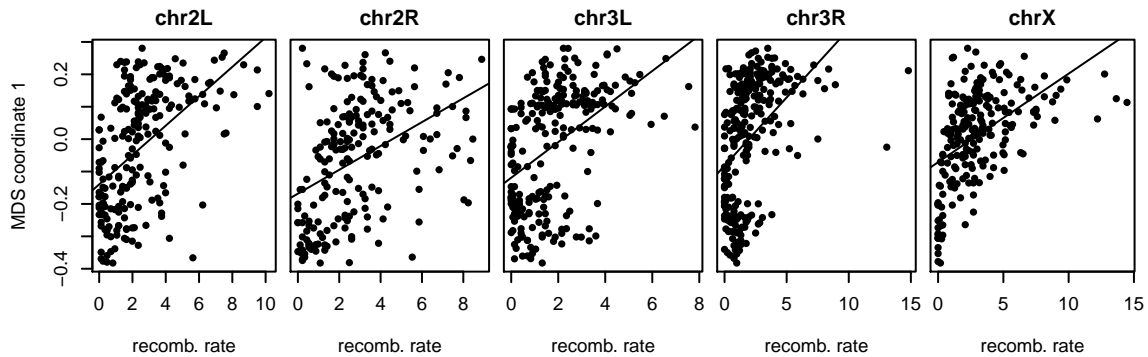


Figure S11: Recombination rate, and the effects of population structure for *Drosophila melanogaster*: this shows the first MDS coordinate and recombination rate (in cM/Mbp), as in Figure 4, against each other. Since the windows underlying estimates of Figure 4 do not coincide, to obtain correlations we divided the genome into 100Kbp bins, and for each variable (recombination rate and MDS coordinate 1) averaged the values of each overlapping bin with weight proportional to the proportion of overlap. The correlation coefficient and  $p$ -values for each linear regression are as follows: 2L: correlation = 0.52,  $r^2 = 0.27$ ; 2R: correlation = 0.43,  $r^2 = 0.18$ ; 3L: correlation = 0.47,  $r^2 = 0.21$ ; 3R: correlation = 0.46,  $r^2 = 0.21$ ; X: correlation = 0.50,  $r^2 = 0.24$ .

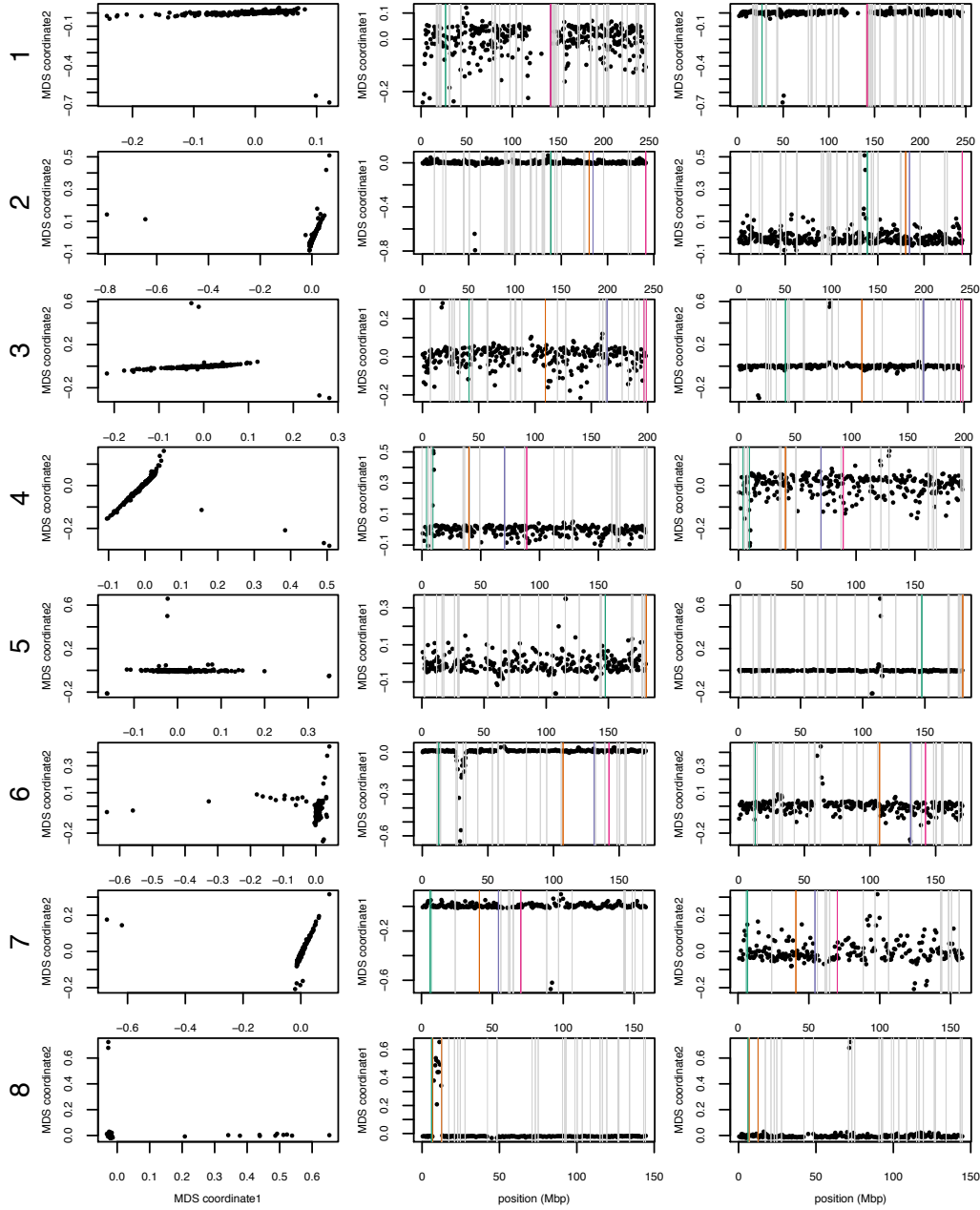


Figure S12: MDS plots for human chromosomes 1-8. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the midpoint of each window against the two MDS coordinates; rows correspond to chromosomes. Colorful vertical lines show the breakpoints of known valid inversions, while grey vertical lines show the breakpoints of predicted inversions.

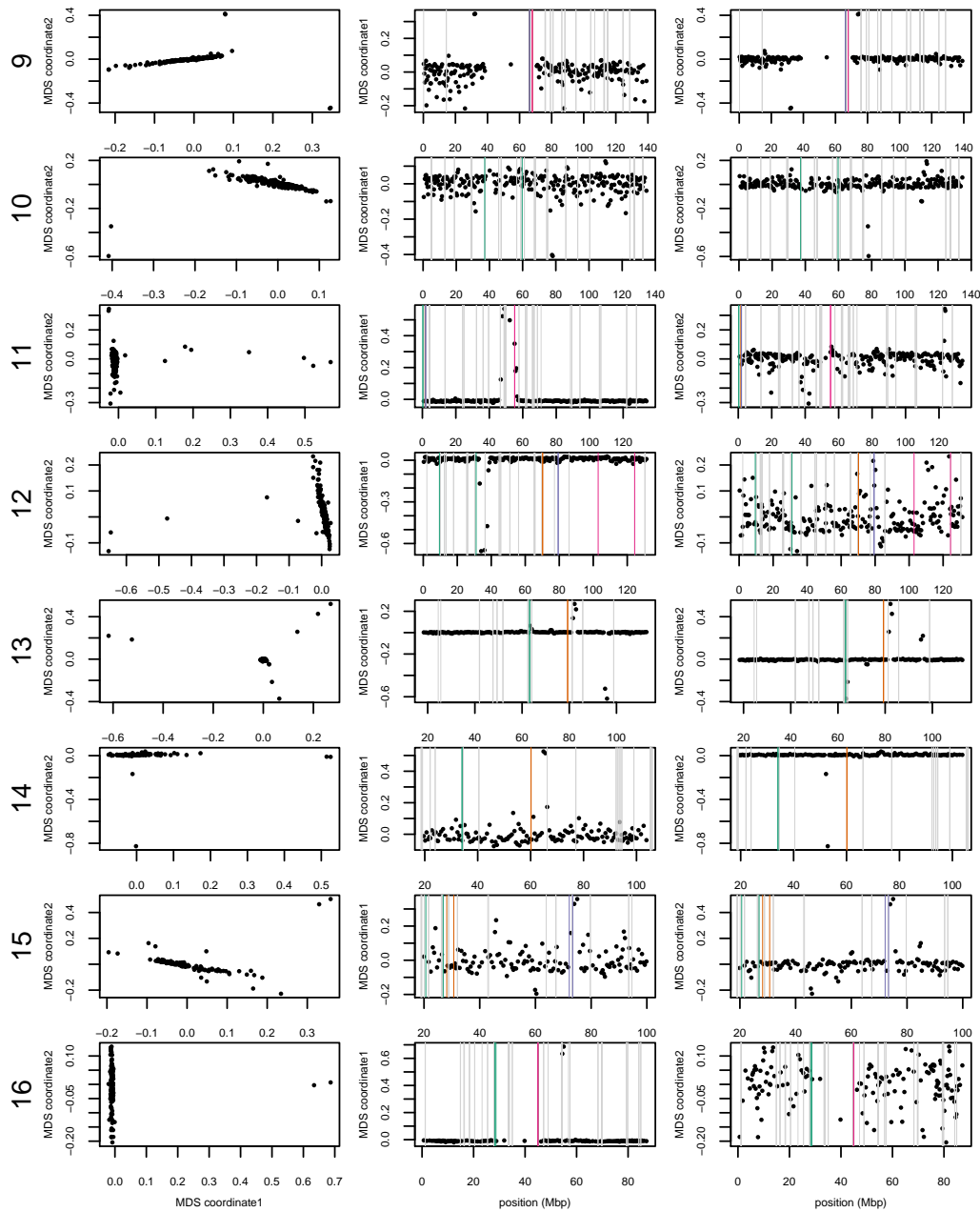


Figure S13: MDS plots for human chromosomes 9-16, as in Supplemental Figure S12.

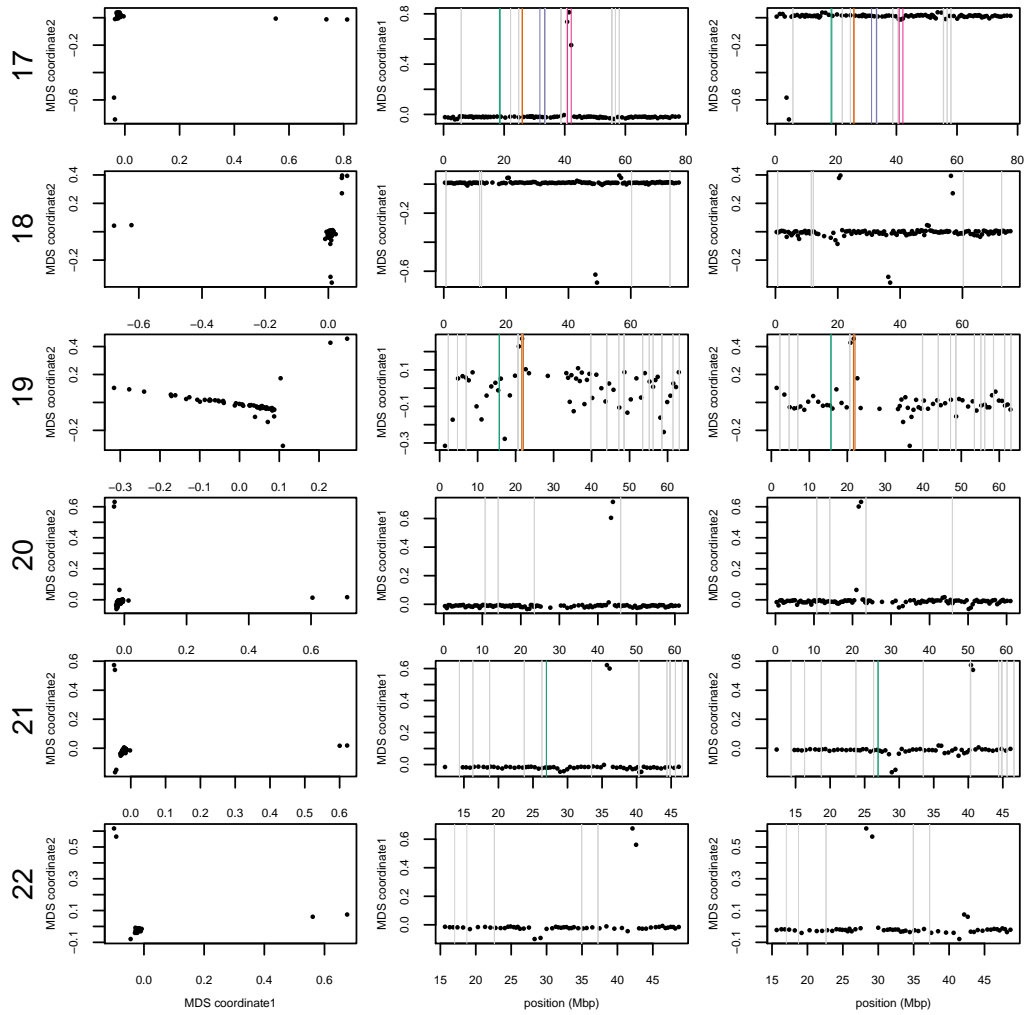


Figure S14: MDS plots for human chromosomes 17-22, as in Supplemental Figure S12.



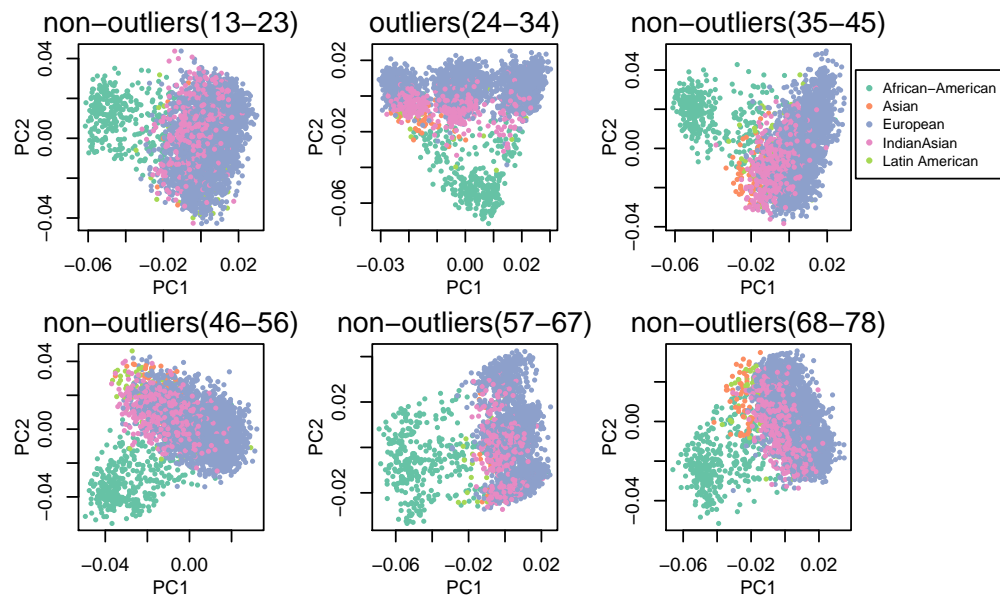


Figure S15: Comparison of PCA figures within outlying windows (center column) and flanking non-outlying windows (left and right columns) for the two windows having outlying MDS scores on chromosome 8.

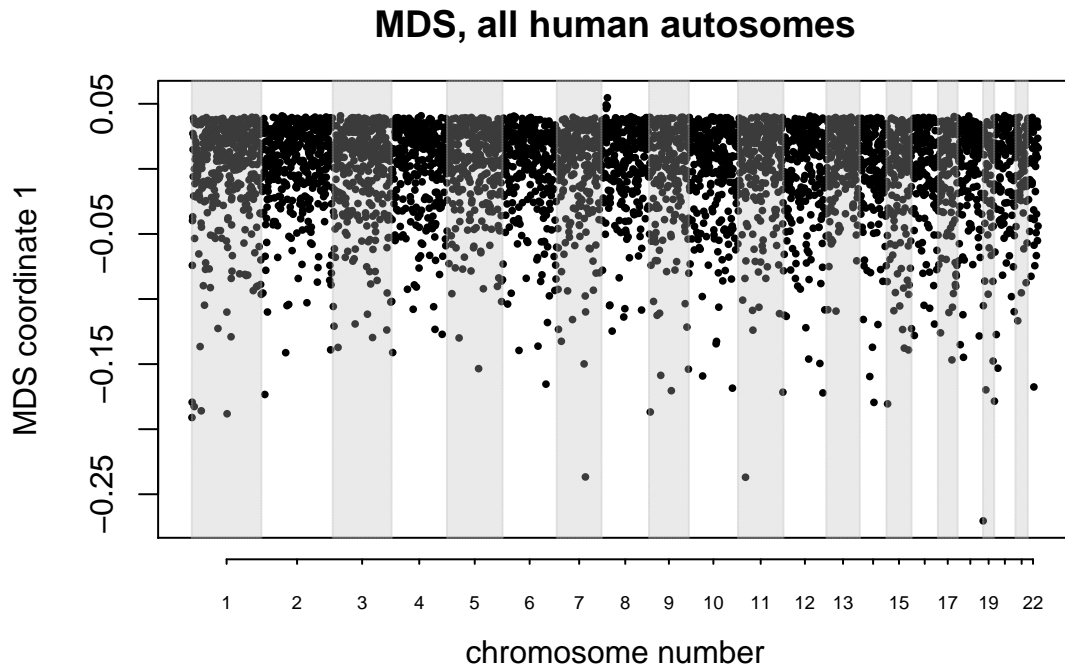


Figure S16: MDS visualization of variation in the effects of population structure amongst windows across *all* human autosomes simultaneously. The small group of windows with positive outlying MDS values lie around the inversion at 8p23.

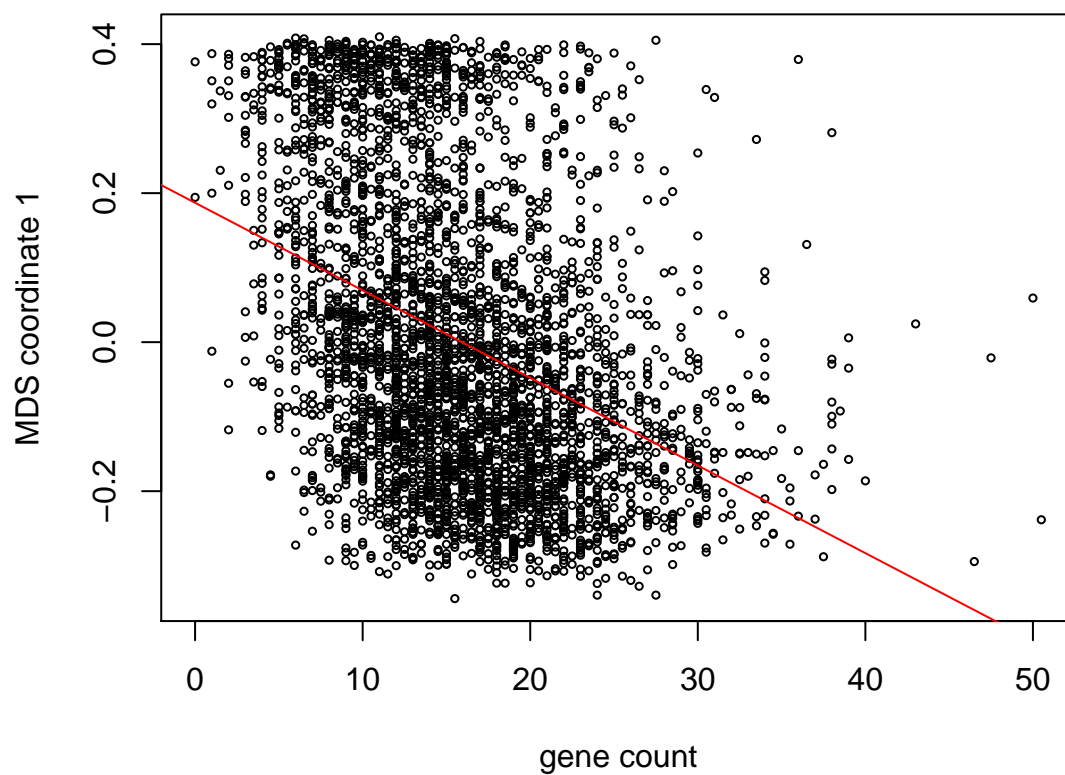


Figure S17: First MDS coordinate against gene density for all 8 chromosomes of *M. truncatula*. The first MDS coordinate is significantly correlated with gene count ( $r = 0.149$ ,  $p = 2.2 \times 10^{-16}$ ).

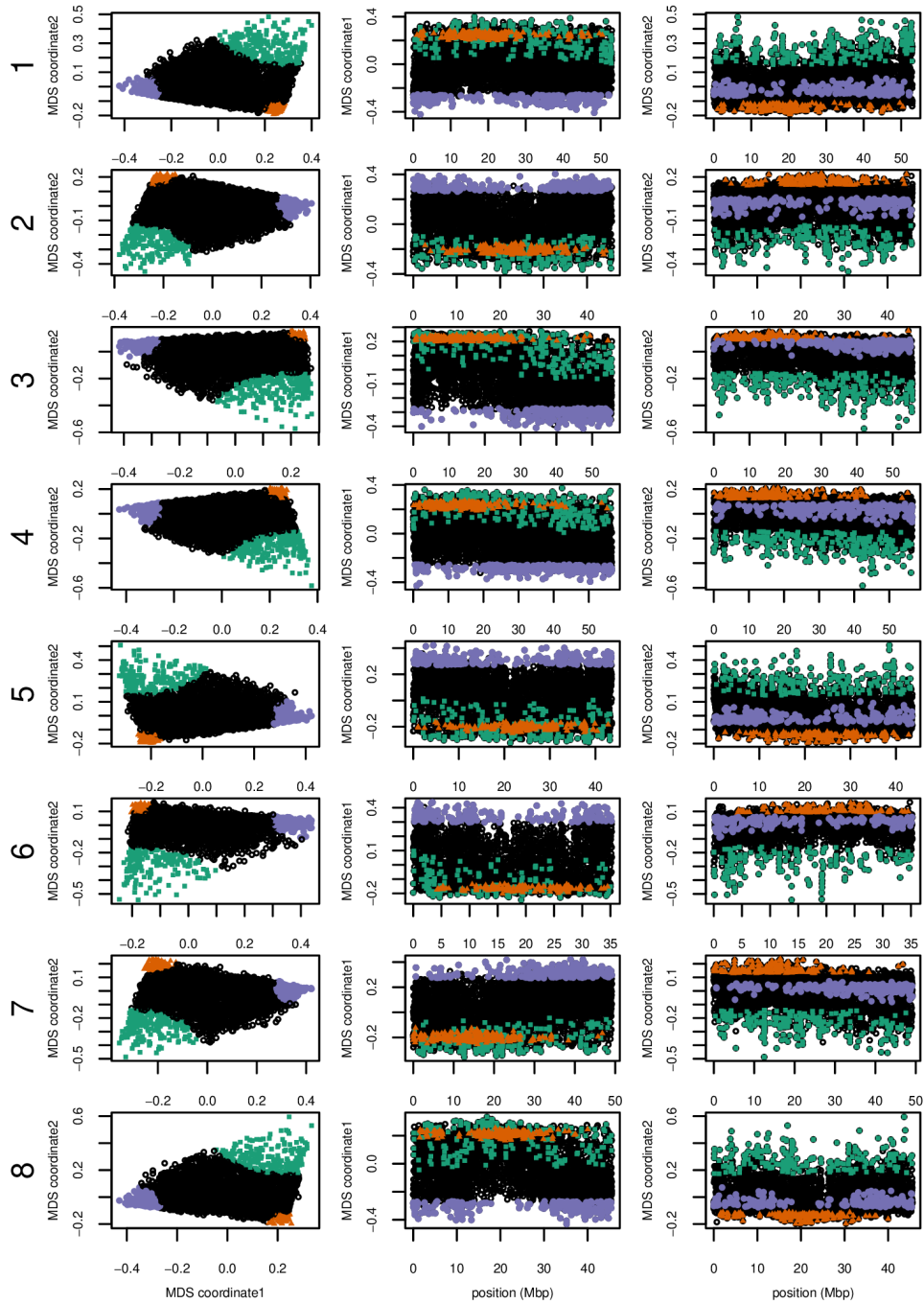


Figure S18: MDS visualizations of the effects of population structure for all 8 chromosomes of the *Medicago truncatula* data, using windows of  $10^4$  SNPs.

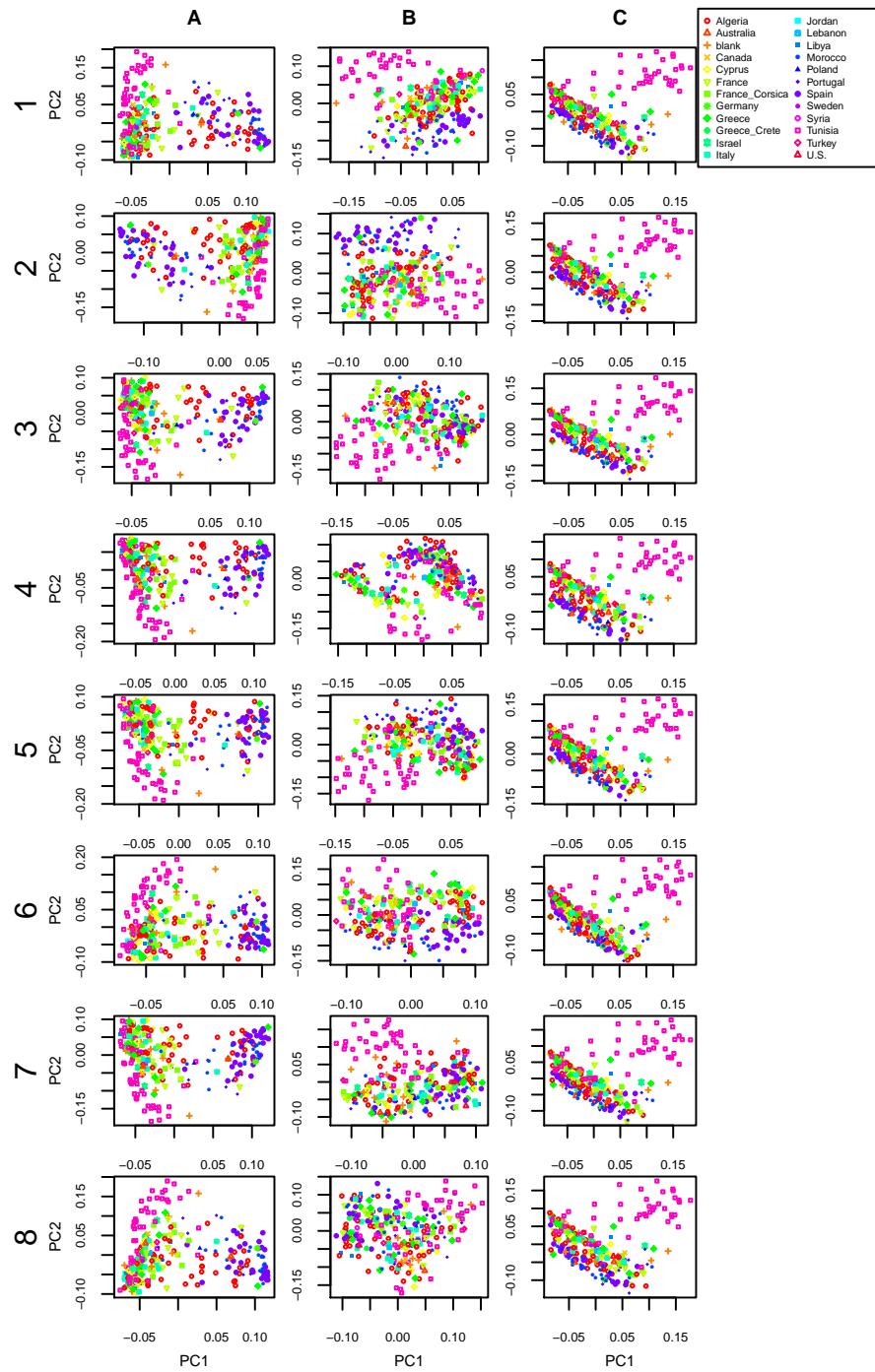


Figure S19: PCA plots for regions colored in Figure S18 on all 8 chromosomes of *Medicago truncatula*: (A) green, (B) orange, and (C) purple.