

Space is the Place: Effects of Continuous Spatial Structure on Analysis of Population Genetic Data

C.J. Battey*,¹, Peter L. Ralph*,[†] and Andrew D. Kern*,[†]

*University of Oregon Dept. Biology, Institute for Ecology Evolution

ABSTRACT Real geography is continuous, but standard models in population genetics are based on discrete, well-mixed populations. As a result many methods of analyzing genetic data assume that samples are a random draw from a well-mixed population, but are applied to clustered samples from populations that are structured clinally over space. Here we use simulations of populations living in continuous geography to study the impacts of dispersal and sampling strategy on population genetic summary statistics, demographic inference, and genome-wide association studies. We find that most common summary statistics have distributions that differ substantially from that seen in well-mixed populations, especially when Wright's neighborhood size is less than 100 and sampling is spatially clustered. The combination of low dispersal and clustered sampling causes demographic inference from the site frequency spectrum to infer more turbulent demographic histories, but averaged results across multiple simulations were surprisingly robust to isolation by distance. We also show that the combination of spatially autocorrelated environments and limited dispersal causes genome-wide association studies to identify spurious signals of genetic association with purely environmentally determined phenotypes, and that this bias is only partially corrected by regressing out principal components of ancestry. Last, we discuss the relevance of our simulation results for inference from genetic variation in real organisms.

KEYWORDS Space; Population Structure; Demography; Haplotype block sharing; GWAS

Introduction

The inescapable reality that biological organisms live, move, and reproduce in continuous geography is usually omitted from population genetic models. However, mates tend to live near to one another and to their offspring, leading to a positive correlation between genetic differentiation and geographic distance. This pattern of "isolation by distance" (Wright 1943) is one of the most widely replicated empirical findings in population genetics (Aguillon *et al.* 2017; Jay *et al.* 2012; Sharbel *et al.* 2000). Despite a long history of analytical work describing the genetics of populations distributed across continuous geography (e.g., Wright (1943); Rousset (1997); Barton *et al.* (2002, 2010); Ringbauer *et al.* (2017); Robledo-Arnuncio and Rousset (2010); Wilkins and Wakeley (2002); Wilkins (2004a)), much modern work still describes geographic structure as a set of discrete populations connected by migration (e.g., Wright 1931; Epperson 2003; Rousset and Leblois 2011; Shirk and Cushman 2014; Lundgren and Ralph 2018; Al-Asadi *et al.* 2019). For this reason, most population genetics statistics are interpreted with reference to discrete, well-mixed populations, and most empirical papers analyze variation within clusters of

Manuscript compiled: Monday 18th November, 2019

¹301 Pacific Hall, University of Oregon Dept. Biology, Institute for Ecology and Evolution. cbattey2@uoregon.edu.

[†]these authors co-supervised this project

40 genetic variation inferred by programs like *STRUCTURE* (Pritchard *et al.* 2000) with methods that
41 assume these are randomly mating units.

42 The assumption that populations are “well-mixed” has important implications for downstream
43 inference of selection and demography. Methods based on the coalescent (Kingman 1982; Wakeley
44 2009) assume that the sampled individuals are a random draw from a well-mixed population that is
45 much larger than the sample (Wakeley and Takahashi 2003). The key assumption is that the individuals
46 of each generation are *exchangeable*, so that there is no correlation between the fate or fecundity of a
47 parent and that of their offspring (Huillet and Möhle 2011). If dispersal or mate selection is limited by
48 geographic proximity, this assumption can be violated in many ways. For instance, if mean viability or
49 fecundity is spatially autocorrelated, then limited geographic dispersal will lead to parent–offspring
50 correlations. Furthermore, nearby individuals will be more closely related than an average random
51 pair, so drawing multiple samples from the same area of the landscape will represent a biased sample
52 of the genetic variation present in the whole population (Städler *et al.* 2009).

53 Two areas in which spatial structure may be particularly important are demographic inference and
54 genome-wide association studies (GWAS). Previous work has found that discrete population structure
55 can create false signatures of population bottlenecks when attempting to infer demographic histories
56 from microsatellite variation (Chikhi *et al.* 2010), statistics summarizing the site frequency spectrum
57 (SFS) (Ptak and Przeworski 2002; Städler *et al.* 2009; St. Onge *et al.* 2012), or runs of homozygosity in a
58 single individual (Mazet *et al.* 2015). The increasing availability of whole-genome data has led to the
59 development of many methods that attempt to infer detailed trajectories of population sizes through
60 time based on a variety of summaries of genetic data (Liu and Fu 2015; Schiffels and Durbin 2014;
61 Sheehan *et al.* 2013; Terhorst *et al.* 2016). Because all of these methods assume that the populations
62 being modeled are approximately randomly mating, they are likely affected by spatial biases in the
63 genealogy of sampled individuals (Wakeley 1999), which may lead to incorrect inference of population
64 changes over time (Mazet *et al.* 2015). However, previous investigations of these effects have focused on
65 discrete rather than continuous space models, and the level of isolation by distance at which inference
66 of population size trajectories become biased by structure is not well known. Here we test how two
67 methods suitable for use with large samples of individuals – stairwayplot (Liu and Fu 2015) and
68 SMC++ (Terhorst *et al.* 2016) – perform when applied to populations evolving in continuous space
69 with varying sampling strategies and levels of dispersal.

70 Spatial structure is also a major challenge for interpreting the results of genome-wide association
71 studies (GWAS). This is because many phenotypes of interest have strong geographic differences due
72 to the (nongenetic) influence of environmental or socioeconomic factors, which can therefore show
73 spurious correlations with spatially patterned allele frequencies (Bulik-Sullivan *et al.* 2015; Mathieson
74 and McVean 2012). Indeed, two recent studies found that previous evidence of polygenic selection on
75 human height in Europe was confounded by subtle population structure (Sohail *et al.* 2018; Berg *et al.*
76 2018), suggesting that existing methods to correct for population structure in GWAS are insufficient.
77 However we have little quantitative idea of the population and environmental parameters that can be
78 expected to lead to biases in GWAS.

79 Last, some of the most basic tools of population genetics are summary statistics like F_{IS} and
80 Tajima’s D , which are often interpreted as reflecting the influence of selection or demography on
81 sampled populations (Tajima 1989). Statistics like Tajima’s D are essentially summaries of the site
82 frequency spectrum, which itself reflects variation in branch lengths and tree structure of the underlying
83 genealogies of sampled individuals. Geographically limited mate choice distorts the distribution of
84 these genealogies (Maruyama 1972; Wakeley 1999), which can affect the value of Tajima’s D (Städler
85 *et al.* 2009). Similarly, the distribution of tract lengths of identity by state among individuals contains
86 information about not only historical demography (Harris and Nielsen 2013; Ralph and Coop 2013)
87 and selection (Garud *et al.* 2015), but also dispersal and mate choice (Ringbauer *et al.* 2017; Baharian
88 *et al.* 2016). We are particularly keen to examine how such summaries will be affected by models that
89 incorporate continuous space, both to evaluate the assumptions underlying existing methods and to
90 identify where the most promising signals of geography lie.

91 To study this, we have implemented an individual-based model in continuous geography that
92 incorporates overlapping generations, local dispersal of offspring, and density-dependent survival. We

93 simulate chromosome-scale genomic data in tens of thousands of individuals from parameter regimes
94 relevant to common subjects of population genetic investigation such as humans and *Drosophila*, and
95 output the full genealogy and recombination history of all final-generation individuals. We use these
96 simulations to test how sampling strategy interacts with geographic population structure to cause
97 systematic variation in population genetic summary statistics typically analyzed assuming discrete
98 population models. We then examine how the fine-scale spatial structures occurring under limited
99 dispersal impact demographic inference from the site frequency spectrum. Last, we examine the
100 impacts of continuous geography on genome-wide association studies (GWAS) and identify regions of
101 parameter space under which the results from GWAS may be misleading.

102 Materials and Methods

103 **Modeling Evolution in Continuous Space**

104 The degree to which genetic relationships are geographically correlated depends on the chance that
105 two geographically nearby individuals are close relatives – in modern terms, by the tension between
106 migration (the chance that one is descended from a distant location) and coalescence (the chance that
107 they share a parent). A key early observation by Wright (Wright 1946) is that this balance is often
108 nicely summarized by the “neighborhood size”, defined to be $N_W = 4\pi\rho\sigma^2$, where σ is the mean
109 parent–offspring distance and ρ is population density. This can be thought of as proportional to the
110 average number of potential mates for an individual (those within distance 2σ), or the number of
111 potential parents of a randomly chosen individual. Empirical estimates of neighborhood size vary
112 hugely across species – even in human populations, estimates range from 40 to over 5,000 depending
113 on the population and method of estimation (Table 1).

114 The first approach to modeling continuously distributed populations was to endow individuals in a
115 Wright-Fisher model with locations in continuous space. However, since the total size of the population
116 is constrained, this introduces interactions between arbitrarily distant individuals, which (aside from
117 being implausible) was shown by Felsenstein (1975) to eventually lead to unrealistic population
118 clumping if the range is sufficiently large. Another method for modeling spatial populations is to
119 assume the existence of a grid of discrete randomly mating populations connected by migration, thus
120 enforcing regular population density by edict. Among many other results drawn from this class of
121 “lattice” or “stepping stone” models (Epperson 2003), Rousset (1997) showed that the slope of the linear
122 regression of genetic differentiation (F_{ST}) against the logarithm of spatial distance is an estimate of
123 neighborhood size. Although these grid models may be good approximations of continuous geography
124 in many situations, they do not model demographic fluctuations, and limit investigation of spatial
125 structure below the level of the deme, assumptions whose impacts are unknown. An alternative
126 method for dealing with continuous geography is a new class of coalescent models, the Spatial Lambda
127 Fleming-Viot models (Barton *et al.* 2010; Kelleher *et al.* 2014).

128 To avoid hard-to-evaluate approximations, we here used forward-time, individual-based simulations.
129 The question of what regulates real populations has a long history and many answers (e.g.,
130 Lloyd 1967; Antonovics and Levin 1980; Crawley 1990), but it is clear that populations must at some
131 point have density-dependent feedback on population size, or else they would face eventual extinction
132 or explosion. In the absence of unrealistic global population regulation, this regulation must be local,
133 and there are many ways to arrange this (Bolker *et al.* 2003). In our simulations, we decrease each
134 individual’s probability of survival with increasing local population density, which shifts reproduc-
135 tive output towards low-density regions (and prevents the unrealistic population clustering seen by
136 Felsenstein (1975), see Supplemental Figure S1). Such models have been used extensively in ecological
137 modeling (Durrett and Levin 1994; Bolker and Pacala 1997; Law *et al.* 2003; Fournier and Méléard 2004;
138 Champer *et al.* 2019) but rarely in population genetics, where to our knowledge implementations of
139 continuous space models before their availability through SLiM (Haller and Messer 2019) have focused
140 on a small number of genetic loci (e.g., Slatkin and Barton 1989; Barton *et al.* 2002; Robledo-Arnuncio
141 and Rousset 2010; Rossine 2014), which limits the ability to investigate the impacts of continuous space
142 on genome-wide genetic variation as is now routinely sampled from real organisms. By simulating
143 chromosome-scale sequence alignments and complete population histories we are able to treat our

144 simulations as real populations and replicate the sampling designs and analyses commonly conducted
145 on real genomic data.

146 **A Forward-Time Model of Evolution in Continuous Space**

147 We simulated populations using the non-Wright-Fisher module in the program SLiM v3.1 (Haller and
148 Messer 2019). Each time step consists of three stages: reproduction, dispersal, and mortality. To reduce
149 the parameter space we use the same parameter, denoted σ , to modulate the spatial scale of interactions
150 at all three stages by adjusting the standard deviation of the corresponding Gaussian functions. As in
151 previous work (Wright 1943; Ringbauer *et al.* 2017), σ is equal to the mean parent-offspring distance.

152 At the beginning of the simulation individuals are distributed uniformly at random on a continuous,
153 square landscape. Individuals are hermaphroditic, and each time step, each produces a Poisson
154 number of offspring with mean $1/L$ where L is the expected lifespan. Offspring disperse a Gaussian-
155 distributed distance away from the parent with mean zero and standard deviation σ in both the x and
156 y coordinates. Each offspring is produced with a mate selected randomly from those within distance
157 3σ , with probability of choosing a neighbor at distance x proportional to $\exp(-x^2/2\sigma^2)$.

158 To maintain a stable population, mortality increases with local population density. To do this we say
159 that individuals at distance d have a competitive interaction with strength $g(d)$, where g is the Gaussian
160 density with mean zero and standard deviation σ . Then, the sum of all competitive interactions with
161 individual i is $n_i = \sum_j g(d_{ij})$, where d_{ij} is the distance between individuals i and j and the sum is over
162 all neighbors within distance 3σ . Since g is a probability density, n_i is an estimate of the number of
163 nearby individuals per unit area. Then, given a per-unit carrying capacity K , the probability of survival
164 until the next time step for individual i is

$$p_i = \min \left(0.95, \frac{1}{1 + n_i / (K(1 + L))} \right). \quad (1)$$

165 We chose this functional form so that the equilibrium population density per unit area is close to K ,
166 and the mean lifetime is around L ; for more description see Appendix .

167 An important step in creating any spatial model is dealing with range edges. Because local popula-
168 tion density is used to model competition, edge or corner populations can be assigned artificially high
169 fitness values because they lack neighbors within their interaction radius but outside the bounds of the
170 simulation. We approximate a decline in habitat suitability near edges by decreasing the probability of
171 survival proportional to the square root of distance to edges in units of σ . The final probability of
172 survival for individual i is then

$$s_i = p_i \min(1, \sqrt{x_i/\sigma}) \min(1, \sqrt{y_i/\sigma}) \min(1, \sqrt{(W - x_i)/\sigma}) \min(1, \sqrt{(W - y_i)/\sigma}) \quad (2)$$

173 where x_i and y_i are the spatial coordinates of individual i , and W is the width (and height) of the
174 square habitat. This buffer roughly counteracts the increase in fitness individuals close to the edge
175 would otherwise have, though the effect is relatively subtle S2.

176 To isolate spatial effects from other components of the model such as overlapping generations,
177 increased variance in reproductive success, and density-dependent fitness, we also implemented
178 simulations identical to those above except that mates are selected uniformly at random from the
179 population, and offspring disperse to a uniform random location on the landscape. We refer to this
180 model as the “random mating” model, in contrast to the first, “spatial” model.

181 We stored the full genealogy and recombination history of final-generation individuals as tree
182 sequences (Kelleher *et al.* 2018), as implemented in SLiM (Haller *et al.* 2019). Scripts for figures and
183 analyses are available at <https://github.com/petrelharp/spaceness>.

184 We ran 400 simulations for the spatial and random-mating models on a square landscape of width
185 $W = 50$ with per-unit carrying capacity $K = 5$ (census $N \approx 10,000$), average lifetime $L = 4$, genome
186 size = 10^8 , recombination rate = 10^{-9} , and drawing σ values from a uniform distribution between 0.2
187 and 4. To speed up the simulations and limit memory overhead we used a mutation rate of 0 in SLiM
188 and later applied mutations to the tree sequence with msprime’s mutate function (Kelleher *et al.* 2016).

189 Because msprime applies mutations proportionally to elapsed time, we divided the mutation rate of
 190 10^{-8} mutations per site per generation by the average generation time estimated for each value of
 191 σ (see ‘Demographic Parameters’ below) to convert the rate to units of mutations per site per unit
 192 time. We show that this procedure produces the same site-frequency spectrum as applying mutations
 193 directly in SLiM in Figure ???. Simulations were run for 1.6 million timesteps (approximately $30N$
 194 generations).

195 To check that our model produces reasonable results, we compared its output to that of a stepping-
 196 stone model implemented in msprime (Kelleher *et al.* 2016). These results are discussed in detail in
 197 Appendix 1, but in general we find that our model produces many of patterns generated by stepping
 198 stone models while avoiding some artifacts associated with discretization of the landscape.

199 Demographic Parameters

200 Our demographic model includes parameters that control population density (K), mean life span (L),
 201 and dispersal distance (σ). However, nonlinearity of local demographic stochasticity causes actual
 202 realized averages of these demographic quantitites to deviate from the specified values in a way
 203 that depends on the neighborhood size. Therefore, to properly compare to theoretical expectations,
 204 we empirically calculated these demographic quantities in simulations. We recorded the census
 205 population size in all simulations. To estimate generation times, we stored ages of the parents of
 206 every new individual born across 200 timesteps, after a 100 generation burn-in, and took the mean. To
 207 estimate variance in offspring number, we tracked the number of offspring for all individuals for 100
 208 timesteps following a 100-timestep burn-in period, subset the resulting table to include only the last
 209 timestep recorded for each individual, and calculated the variance in number of offspring across all
 210 individuals in timesteps 50-100. All calculations were performed with information recorded in the tree
 211 sequence, using pyslim (<https://github.com/tskit-dev/pyslim>).

212 Sampling

213 Our model records the genealogy and sequence variation of the complete population, but in real data,
 214 genotypes are only observed from a relatively small number of sampled individuals. We modeled three
 215 sampling strategies similar to common data collection methods in empirical genetic studies (Figure 1).
 216 “Random” sampling selects individuals at random from across the full landscape, “point” sampling
 217 selects individuals proportional to their distance from four equally spaced points on the landscape,
 218 and “midpoint” sampling selects individuals in proportion to their distance from the middle of the
 219 landscape. Downstream analyses were repeated across all sampling strategies.

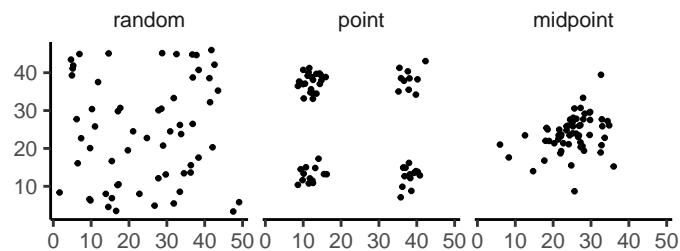


Figure 1 Example sampling maps for 60 individuals on a 50×50 landscape for midpoint, point, and random sampling strategies, respectively.

220 Summary Statistics

221 We calculated the site frequency spectrum and a set of 18 summary statistics (Table S1) from 60 diploid
 222 individuals sampled from the final generation of each simulation using the python package scikit-allel
 223 (Miles and Harding 2017). Statistics included common single-population summaries including mean
 224 pairwise divergence (π), inbreeding coefficient (F_{IS}), and Tajima’s D , as well as an isolation-by-distance

225 regression of genetic distance (D_{xy}) against the logarithm of geographic distance analogous to Rousset
226 (1997)'s approach, which we summarized as the correlation coefficient between the logarithm of the
227 spatial distance and the proportion of identical base pairs across pairs of individuals.

228 Following recent studies that showed strong signals for dispersal and demography in the distribution
229 of shared haplotype block lengths (Ringbauer *et al.* 2017; Baharian *et al.* 2016), we also calculated
230 various summaries of the distribution of pairwise identical-by-state (IBS) block lengths among sam-
231 pled chromosomes. The full distribution of lengths of IBS tracts for each pair of chromosomes was
232 first calculated with a custom python function. We then calculated the first three moments of this
233 distribution (mean, variance, and skew) and the number of blocks over 10^6 base pairs both for each
234 pair of individuals and for the full distribution across all pairwise comparisons.

235 We then estimated correlation coefficients between spatial distance and each moment of the pairwise
236 IBS tract distribution. Because more closely related individuals on average share longer haplotype
237 blocks we expect that spatial distance will be negatively correlated with mean haplotype block length,
238 and that this correlation will be strongest (i.e., most negative) when dispersal is low. The variance,
239 skew, and count of long haplotype block statistics are meant to reflect the relative length of the right
240 (upper) tail of the distribution, which represents the frequency of long haplotype blocks, and so should
241 reflect recent demographic events (Chapman and Thompson 2002). For a subset of simulations, we
242 also calculated cumulative distributions for IBS tract lengths across pairs of distant (> 48 map units)
243 and nearby (< 2 map units) individuals. Last, we examined the relationship between allele frequency
244 and the spatial dispersion of an allele by calculating the average distance among individuals carrying
245 each derived allele in a set of simulations representing a range of neighborhood sizes.

246 The effects of sampling on summary statistic estimates were summarized by testing for differences
247 in mean (ANOVA, (R Core Team 2018)) and variance (Levene's test, (Fox and Weisberg 2011)) across
248 sampling strategies for each summary statistic.

249 Last, we used a subset of summary statistics to compare how our continuous model performs
250 relative to a reverse-time stepping stone model (Appendix 1).

251 **Demographic Inference**

252 To assess the impacts of continuous spatial structure on demographic inference we inferred pop-
253 ulation size histories for all simulations using two approaches: stairwayplot (Liu and Fu 2015) and
254 SMC++ (Terhorst *et al.* 2016). Stairwayplot fits its model to a genome-wide estimate of the SFS, while
255 SMC++ also incorporates linkage information. For both methods we sampled 20 individuals from all
256 spatial simulations using random, midpoint, and point sampling strategies.

257 As recommended by its documentation, we used stairwayplot to fit models with multiple bootstrap
258 replicates drawn from empirical genomic data, and took the median inferred N_e per unit time as
259 the best estimate. We calculated site frequency spectra with scikit-allel (Miles and Harding 2017),
260 generated 100 bootstrap replicates per simulation by resampling over sites, and fit models for all
261 bootstrap samples using default settings.

262 For SMC++, we first output genotypes as VCF with msprime and then used SMC++'s standard
263 pipeline for preparing input files assuming no polarization error in the SFS. We used the first individual
264 in the VCF as the "designated individual" when fitting models, and allowed the program to estimate
265 the recombination rate during optimization. We fit models using the 'estimate' command rather than
266 the now recommended cross-validation approach because our simulations had only a single contig.

267 To evaluate the performance of these methods we binned simulations by neighborhood size, took a
268 rolling median of inferred N_e trajectories across all model fits in a bin for each method and sampling
269 strategy. We also examined how varying levels of isolation by distance impacted the variance of N_e
270 estimates by calculating the standard deviation of N_e from each best-fit model and plotting these
271 against neighborhood size.

273 **Association Studies**

274 To assess the degree to which spatial structure confounds GWAS we simulated four types of nongenetic
275 phenotype variation for 1000 randomly sampled individuals in each spatial SLiM simulation and

276 conducted a linear regression GWAS with principal components as covariates in PLINK (Purcell *et al.* 277 2007). SNPs with a minor allele frequency less than 0.5% were excluded from this analysis. Phenotype 278 values were set to vary by two standard deviations across the landscape in a rough approximation 279 of the variation seen in height across Europe (Turchin *et al.* 2012; Garcia and Quintana-Domeque 280 2006, 2007). Conceptually our approach is similar to that taken by Mathieson and McVean (2012), 281 though here we model fully continuous spatial variation and compare GWAS output across a range of 282 dispersal distances.

283 In all simulations, the phenotype of each individual is determined by adding independent Gaussian 284 noise with mean zero to a mean that may depend on spatial position such that the mean of spatially- 285 correlated phenotypes varies by two standard deviations across the landscape. We then adjust the 286 geographic pattern of mean phenotype to create four spatially autocorrelated environmental influences 287 on phenotype. In the first simulation of *nonspatial* environments, the mean did not change, so that all 288 individuals' phenotypes were drawn independently from a Gaussian distribution with mean 110 and 289 standard deviation 10. Next, to simulate *clinal* environmental influences on phenotype, we increased 290 the mean phenotype from 100 on the left edge of the range to 120 on the right edge (two phenotypic 291 standard deviations). Concretely, an individual at position (x, y) in a 50×50 landscape has mean 292 phenotype $100 + 2x/5$. Third, we simulated a more concentrated "*corner*" environmental effect by 293 setting the mean phenotype for individuals with both x and y coordinates below 20 to 120 (two 294 standard deviations above the rest of the map). Finally, in "*patchy*" simulations we selected 10 random 295 points on the map and set the mean phenotype of all individuals within three map units of each of 296 these points to 120.

297 We performed principal components analysis (PCA) using scikit-allel (Miles and Harding 2017) on 298 the matrix of derived allele counts by individual for each simulation. SNPs were first filtered to remove 299 strongly linked sites by calculating LD between all pairs of SNPs in a 200-SNP moving window and 300 dropping one of each pair of sites with an R^2 over 0.1. The LD-pruned allele count matrix was then 301 centered and all sites scaled to unit variance when conducting the PCA, following recommendations 302 in Patterson *et al.* (2006).

303 We ran linear-model GWAS both with and without the first 10 principal components as covariates 304 in PLINK and summarized results across simulations by counting the number of SNPs with p -value 305 below 0.05 after adjusting for an expected false positive rate of less than 5% (Benjamini and Yekutieli 306 2001). We also examined p values for systematic inflation by estimating the expected values from a 307 uniform distribution (because no SNPs were used when generating phenotypes), plotting observed 308 against expected values for all simulations, and summarizing across simulations by finding the mean 309 σ value in each region of quantile-quantile space. Results from all analyses were summarized and 310 plotted with the "ggplot2" (Wickham 2016) and "cowplot" (Wilke 2019) packages in R (R Core Team 311 2018).

312 **Results**

313 **Demographic Parameters and Run Times**

314 Adjusting the spatial dispersal and interaction distance, σ , has a surprisingly large effect on de- 315 mographic quantities that are usually fixed in Wright-Fisher models – the generation time, census 316 population size, and variance in offspring number, shown in Figure 2. Because our simulation is 317 parameterized on an individual level, these population parameters emerge as a property of the interac- 318 tions among individuals rather than being directly set. Variation across runs occurs because, even 319 though the parameters K and L that control population density and mean lifetime respectively were 320 the same in all simulations, the strength of stochastic effects depends strongly on the spatial interaction 321 distance σ . For instance, the population density near to individual i (denoted n_i above) is computed 322 by averaging over roughly $N_W = 4\pi K\sigma^2$ individuals, and so has standard deviation proportional to 323 $1/\sqrt{N_W}$ – it is more variable at lower densities. (Recall that N_W is Wright's neighborhood size.) Since 324 the probability of survival is a nonlinear function of n_i , actual equilibrium densities and lifetimes differ 325 from K and L . This is the reason that we included *random mating* simulations – where mate choice and 326 offspring dispersal are both nonspatial – since this should preserve the random fluctuations in local

327 population density while destroying any spatial genetic structure. We verified that random mating
 328 models retained no geographic signal by showing that summary statistics did not differ significantly
 329 between sampling regimes (Table S2), unlike in spatial models (discussed below).

330 There are a few additional things to note about Figure 2. First, all three quantities are non-monotone
 331 with neighborhood size. Census size largely declines as neighborhood size increases for both the spatial
 332 and random mating models. However, for spatial models this decline only begins for neighborhood
 333 size ≥ 10 . By a neighborhood sizes larger than 100, the spatial and random mating models are
 334 indistinguishable from one another, a sign that our simulations are performing as expected. Census
 335 sizes range from $\approx 14,000$ at low σ in the random mating model to $\approx 10,000$ for both models when
 336 neighborhood sizes approach 1,000. The scaling of census sizes in both random-mating and spatial
 337 models appears to be related to two consequences of the spatial competition function: the decline
 338 of fitness at range edges, which effectively reduces the habitable area by one σ around the edge of
 339 the map and so results in a smaller habitable area at high σ values; and variation in the equilibrium
 340 population density given varying competition radii. For spatial models we see an additional effect in
 341 which census size initially rises as neighborhood sizes scale from 2 - 10. This may reflect an Allee effect
 342 (Allee *et al.* 1949) in which some individuals are unable to find mates when the mate selection radius is
 343 very small.

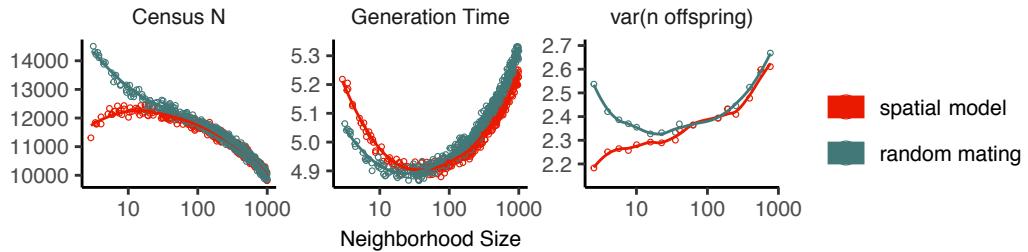


Figure 2 Genealogical parameters from spatial and random mating SLiM simulations, by neighborhood size.

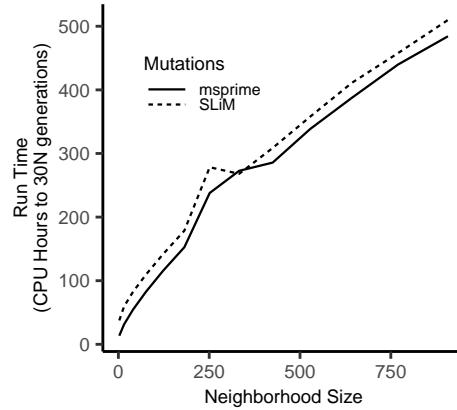


Figure 3 Run times of continuous space simulations with landscape width 50 and expected density 5 under varying neighborhood size. Times are shown for simulations run with mutations applied directly in SLiM (dashed lines) or later applied to tree sequences with msprime (solid lines). Times for simulations run with tree sequence recording disabled are shown in grey.

344 Generation time similarly shows complex behavior with respect to neighborhood sizes, and varies
 345 between 5.2 and 4.9 timesteps per generation across the parameter range explored. Under both the

346 spatial and random mating models, generation time reaches a minimum at a neighborhood size of
347 around 50. Interestingly, under the range of neighborhood sizes that we examined, generation times
348 between the random mating and spatial models are never quite equivalent – presumably this would
349 cease to be the case at neighborhood sizes higher than we simulated here.

350 Last, we looked at the variance in number of offspring – a key parameter determining the effective
351 population size. Surprisingly, the spatial and random mating models behave quite differently: while
352 the variance in offspring number increases nearly monotonically under the spatial model, the random
353 mating model actually shows a decline in the variance in offspring number until a neighborhood size
354 ≈ 10 before it increases and eventually equals what we observe in the spatial case.

355 Run times for our model scale approximately linearly with neighborhood size (Figure 3), with the
356 lowest neighborhood sizes reaching 30N generations in around a day and those with neighborhood
357 size approaching 1,000 requiring up to three weeks of computation. As currently implemented running
358 simulations at neighborhood sizes > 1000 to coalescence is likely impractical, though running these
359 models for more limited timescales and then "completing" the simulation by starting a reverse-time
360 simulation from the resulting tree sequence in msprime is possible.

361 **Impacts of Continuous Space on Population Genetic Summary Statistics**

362 Even though certain aspects of population demography depend on the scale of spatial interactions, it
363 still could be that population genetic variation is well-described by a well-mixed population model.
364 Indeed, mathematical results suggest that genetic variation in some spatial models should be well-
365 approximated by a Wright-Fisher population if neighborhood size is large and all samples are geo-
366 graphically widely separated (Wilkins 2004b; Zähle *et al.* 2005). However, the behavior of most
367 common population genetic summary statistics other than Tajima's D (Städler *et al.* 2009) has not yet
368 been described in realistic geographic models. Moreover, as we will show, spatial sampling strategies
369 can affect summaries of genetic variation at least as strongly as the underlying population dynamics.

370 **Site Frequency Spectra and Summaries of Diversity** Figure 4 shows the effect of varying neighbor-
371 hood size and sampling strategy on the site frequency spectrum (Figure 4A) and several standard
372 population genetic summary statistics (Figure 4B). Consistent with findings in island and stepping
373 stone simulations (Städler *et al.* 2009), the SFS shows a significant enrichment of intermediate frequency
374 variants in comparison to the nonspatial expectation. This bias is most pronounced below neighbor-
375 hood sizes ≤ 100 and is exacerbated by midpoint and point sampling of individuals (depicted in
376 Figure 1). Reflecting this, Tajima's D is quite positive in the same situations (Figure 4B). Notably, the
377 point at which Tajima's D approaches 0 differs strongly across sampling strategies – varying from a
378 neighborhood size of roughly 50 for random sampling to at least 1000 for midpoint sampling.

379 One of the most commonly used summaries of variation is Tajima's summary of nucleotide diver-
380 gence, θ_π , calculated as the mean density of nucleotide differences averaged across pairs of samples.
381 As can be seen in Figure 4B, θ_π in the spatial model is inflated by up to three-fold relative to the
382 random mating model. This pattern is opposite the expectation from census population size (Figure 2),
383 because the spatial model has *lower* census size than the random mating model at neighborhood sizes
384 less than 100. Differences between these models likely occur because θ_π is a measure of mean time to
385 most recent common ancestor between two samples, and at small values of σ , the time for dispersal to
386 mix ancestry across the range exceeds the mean coalescent time under random mating. (For instance,
387 at the smallest value of $\sigma = 0.2$, the range is 250 dispersal distances wide, and since the location
388 of a diffusively moving lineage after k generations has variance $k\sigma^2$, it takes around $250^2 = 62500$
389 generations to mix across the range, which is roughly ten times larger than the random mating effective
390 population size). θ_π using each sampling strategy approaches the random mating expectation at its
391 own rate, but by a neighborhood size of around 100 all models are roughly equivalent. Interestingly,
392 the effect of sampling strategy is reversed relative to that observed in Tajima's D – midpoint sampling
393 reaches random mating expectations around neighborhood size 50, while random sampling is inflated
394 until around neighborhood size 100.

395 Values of observed heterozygosity and its derivative F_{IS} also depend heavily on neighborhood size
396 under spatial models as well as the sampling scheme. F_{IS} is inflated above the expectation across



Figure 4 Log-scaled site frequency spectrum (A) and summary statistic distributions (B) by sampling strategy and neighborhood size.

most of the parameter space examined and across all sampling strategies. This effect is caused by a deficit of heterozygous individuals in low-dispersal simulations – a continuous-space version of the Wahlund effect (Wahlund 1928). Indeed, for random sampling under the spatial model, F_{IS} does not approach the random mating equivalent until neighborhood sizes of nearly 1000. On the other hand, the dependency of raw observed heterozygosity on neighborhood size is not monotone. Under midpoint sampling observed heterozygosity is inflated even over the random mating expectation, as a result of the a higher proportion of heterozygotes occurring in the middle of the landscape (Figure S6). This echoes a report from Shirk and Cushman (2014) who observed a similar excess of heterozygosity in the middle of the landscape when simulating under a lattice model.

IBS tracts and correlations with geographic distance We next turn our attention to the effect of geographic distance on haplotype block length sharing, summarized for sets of nearby and distant individuals in Figure 5. There are two main patterns to note. First, nearby individuals share more long IBS tracts than distant individuals (as expected because they are on average more closely related). Second, the difference in the number of long IBS tracts between nearby and distant individuals decreases as neighborhood size increases. This reflects the faster spatial mixing of populations with higher dispersal, which breaks down the correlation between the IBS tract length distribution and geographic distance. This can also be seen in the bottom row of Figure 4B, where the correlation coefficients between the summaries of the IBS tract length distribution (the mean, skew, and count of tracts over 10^6 bp) and geographic distance approaches 0 as neighborhood size increases.

The patterns observed for correlations of IBS tract lengths with geographic distance are similar to those observed in the more familiar regression of allele frequency measures such as D_{xy} (i.e., “genetic distance”) or F_{ST} against geographic distance (Rousset 1997). D_{xy} is positively correlated with the geographic distance between the individuals, and the strength of this correlation declines as dispersal increases (Figure 4B), as expected (Wright 1943; Rousset 1997). This relationship is very similar across random and point sampling strategies, but is weaker for midpoint sampling, perhaps due to a dearth of long-distance comparisons. In much of empirical population genetics a regression of genetic differentiation against spatial distance is a de-facto metric of the significance of isolation by distance. The similar behavior of moments of the pairwise distribution of IBS tract lengths shows why haplotype block sharing has recently emerged as a promising source of information on spatial demography through methods described in Ringbauer *et al.* (2017) and Baharian *et al.* (2016).

Spatial distribution of allele copies Mutations occur in individuals and spread geographically over time. Because low frequency alleles generally represent recent mutations (Sawyer 1977; Griffiths *et al.* 1999), the geographic dispersion of an allele may covary along with its frequency in the population. To visualize this relationship we calculated the average distance among individuals carrying a focal derived allele across simulations with varying neighborhood sizes, shown in Figure 6. On average we find that low frequency alleles are the most geographically restricted, and that the extent to which geography and allele frequency are related depends on the amount of dispersal in the population. For populations with large neighborhood sizes we found that even very low frequency alleles can be found across the full landscape, whereas in populations with low neighborhood sizes the relationship between distance among allele copies and their frequency is quite strong. This is the basic process underlying Novembre and Slatkin’s (2009) method for estimating dispersal distances based on the distribution of low frequency alleles, and also generates the greater degree of bias in GWAS effect sizes for low frequency alleles identified in Mathieson and McVean (2012).

Effects of Space on Demographic Inference

One of the most important uses for population genetic data is inferring demographic history of populations. As demonstrated above, the site frequency spectrum and the distribution of IBS tracts varies across neighborhood sizes and sampling strategies. Does this variation lead to different inferences of past population sizes? To ask this we inferred population size histories from samples drawn from our simulated populations with two approaches: stairwayplot (Liu and Fu 2015), which uses a genome-wide estimate of the SFS, and SMC++ (Terhorst *et al.* 2016), which incorporates information on both the

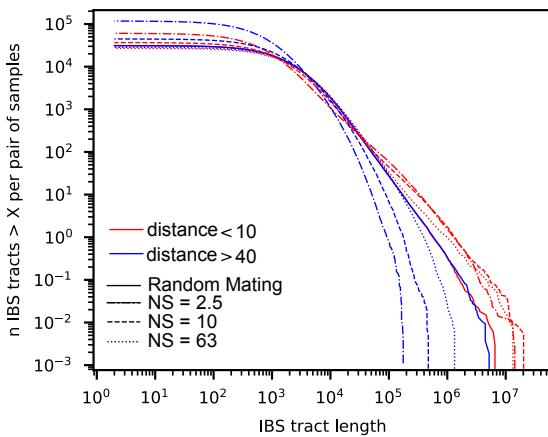


Figure 5 Cumulative distributions for IBS tract lengths per pair of individuals at different geographic distances, across three neighborhood sizes (NS).

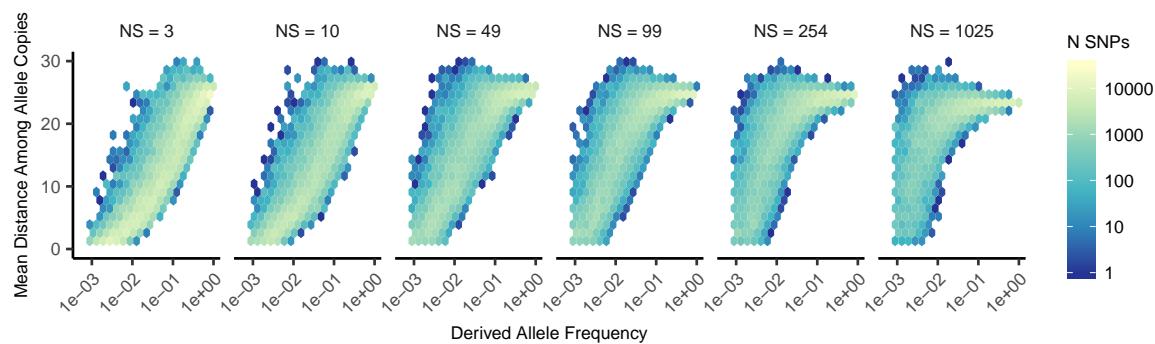


Figure 6 Trends in the distance among allele copies at varying derived allele frequencies and neighborhood sizes (NS).

447 SFS and linkage disequilibrium across the genome.

448 Figure 7A shows the median inferred population size histories from each method across all sim-
449 ulations, grouped by neighborhood size and sampling strategy. In general these methods tend to
450 slightly overestimate ancient population sizes and infer recent population declines when neighborhood
451 sizes are below 20 and sampling is spatially clustered (Figure 7A, Figure S7). The overestimation
452 of ancient population sizes however is relatively minor, averaging around a two-fold inflation at
453 10,000 generations before present in the worst-affected bins. For stairwayplot we found that many
454 runs infer dramatic population bottlenecks in the last 1,000 generations when sampling is spatially
455 concentrated, resulting in ten-fold or greater underestimates of recent population sizes. However
456 SMC++ appeared more robust to this error, with runs on point- and midpoint-sampled simulations at
457 the lowest neighborhood sizes underestimating recent population sizes by roughly half and those on
458 randomly sampled simulations showing little error. Above neighborhood sizes of around 100, both
459 methods performed relatively well when averaging across results from multiple simulations.

460 However, individual model fits from both methods frequently reflected turbulent demographic
461 histories (Figure S7), with the standard deviation of inferred N_e across time points often exceeding
462 the expected N_e for both methods (Figure 7B). That is, despite the constant population sizes in our
463 simulations, both methods tended to infer large fluctuations in population size over time, which could
464 potentially result in incorrect biological interpretations. On average the variance of inferred population
465 sizes was elevated at the lowest neighborhood sizes and declines as dispersal increases, with the
466 strongest effects seen in stairwayplot model fits with for clustered sampling and neighborhood sizes
467 less than 20 (Figure 7B).

468 **GWAS**

469 To ask what confounding effects spatial genetic variation might have on genome-wide association
470 studies we performed GWAS on our simulations using phenotypes that were determined solely by
471 the environment – so, any SNP showing statistically significant correlation with phenotype is a false
472 positive. As expected, spatial autocorrelation in the environment causes spurious associations across
473 much of the genome if no correction for genetic relatedness among samples is performed (Figures 8 and
474 S8). This effect is particularly strong for clinal and corner environments, for which the lowest dispersal
475 levels cause over 60% of SNPs in the sample to return significant associations. Patchy environmental
476 distributions, which are less strongly spatially correlated (Figure 8A), cause fewer false positives
477 overall but still produce spurious associations at roughly 10% of sites at the lowest neighborhood
478 sizes. Interestingly we also observed a small number of false positives in roughly 3% of analyses
479 on simulations with nonspatial environments, both with and without PC covariates included in the
480 regression.

481 The confounding effects of geographic structure are well known, and it is common practice to
482 control for this by including principal components (PCs) as covariates to control for these effects. This
483 mostly works in our simulations – after incorporating the first ten PC axes as covariates, the vast
484 majority of SNPs no longer surpass a significance threshold chosen to have a 5% false discovery rate
485 (FDR). However, a substantial number of SNPs – up to 1.5% at the lowest dispersal distances – still
486 surpass this threshold (and thus would be false positives in a GWAS), especially under “corner” and
487 “patchy” environmental distributions (Figure 8C). At neighborhood sizes larger than 500, up to 0.31%
488 of SNPs were significant for corner and clinal environments. Given an average of 132,000 SNPs across
489 simulations after MAF filtering, this translates to up to 382 false-positive associations; for human-sized
490 genomes, this number would be much larger. In most cases the p values for these associations were
491 significant after FDR correction but would not pass the threshold for significance under the more
492 conservative Bonferroni correction (see example Manhattan plots in figure S8).

493 Clinal environments cause an interesting pattern in false positives after PC correction: at low
494 neighborhood sizes the correction removes nearly all significant associations, but at neighborhood
495 sizes above roughly 250 the proportion of significant SNPs increases to up to 0.4% (Figure 8). This
496 may be due to a loss of descriptive power of the PCs – as neighborhood size increases, the total
497 proportion of variance explained by the first 10 PC axes declines from roughly 10% to 4% (Figure
498 8B). Essentially, PCA seems unable to effectively summarize the weak population structure present in

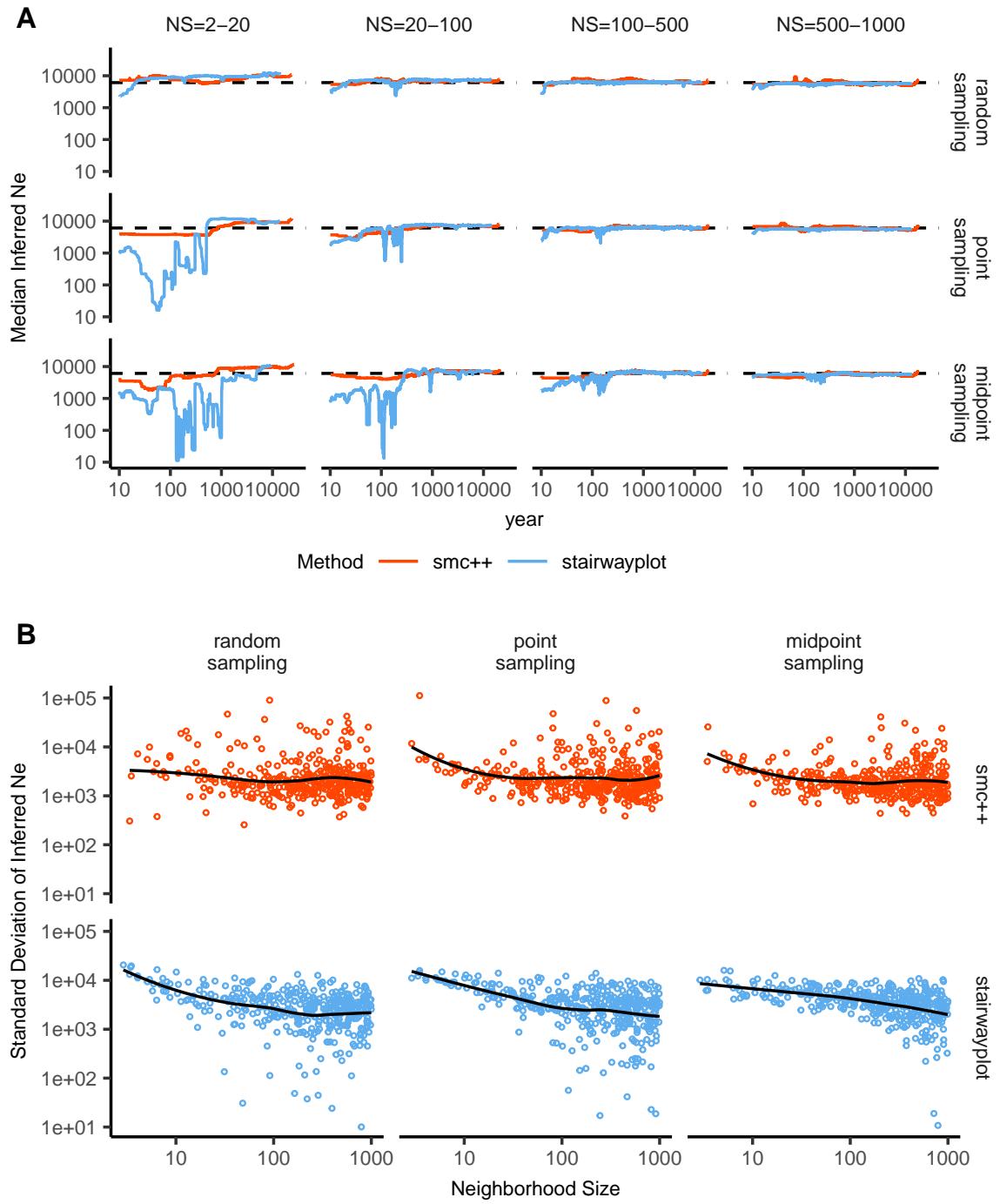


Figure 7 A: Rolling median inferred N_e trajectories for stairwayplot and smc++ across sampling strategies and neighborhood size bins. The dotted line shows the mean N_e of random-mating simulations. B: Standard deviation of individual inferred N_e trajectories, by neighborhood size and sampling strategy. Black lines are loess curves. Plots including individual model fits are shown in Figure S7.

499 large-neighborhood simulations, but these populations continue to have enough spatial structure to
500 create significant correlations between genotypes and the environment. A similar process can also be
501 seen in the corner phenotype distribution, in which the count of significant SNPs initially declines as
502 neighborhood size increases and then increases at approximately the point at which the proportion of
503 variance explained by PCA approaches its minimum.

504 Figure 8D shows quantile-quantile plots that show the degree of genome-wide inflation of test
505 statistics in PC-corrected GWAS across all simulations and environmental distributions. For clinal
506 environments, $-\log_{10}(p)$ values are most inflated when neighborhood sizes are large, consistent with
507 the pattern observed in the count of significant associations after PC regression. In contrast corner
508 and patchy environments cause the greatest inflation in $-\log_{10}(p)$ at neighborhood sizes less than
509 100, which likely reflects the inability of PCA to account for fine-scale structure caused by very limited
510 dispersal. Finally, we observed that PC regression appears to overfit to some degree for all phenotype
511 distributions, visible in Figure 8D as points falling below the 1:1 line.

512 Discussion

513 In this study, we have used efficient forward time population genetic simulations to describe the
514 myriad influence of continuous geography on genetic variation. In particular, we examine how three
515 main types of downstream empirical inference are affected by unmodeled spatial population structure
516 – 1) population genetic summary statistics, 2) inference of population size history, and 3) genome-wide
517 association studies (GWAS). As discussed above, space often matters (and sometimes dramatically),
518 both because of how samples are arranged in space, and because of the inherent patterns of relatedness
519 established by geography.

520 Effects of Dispersal

521 Limited dispersal inflates effective population size, creates correlations between genetic and spatial
522 distances, and introduces strong distortions in the site frequency spectrum that are reflected in a
523 positive Tajima's D (Figure 4). At the lowest dispersal distances, this can increase genetic diversity
524 threefold relative to random-mating expectations. These effects are strongest when neighborhood
525 sizes are below 100, but in combination with the effects of nonrandom sampling they can persist up to
526 neighborhood sizes of at least 1000 (e.g., inflation in Tajima's D and observed heterozygosity under
527 midpoint sampling). If samples are chosen uniformly from across space, the general pattern is similar
528 to expectations of the original analytic model of Wright (1943), which predicts that populations with
529 neighborhood sizes under 100 will differ substantially from random mating, while those above 10,000
530 will be nearly indistinguishable from panmixia.

531 The patterns observed in sequence data reflect the effects of space on the underlying genealogy.
532 Nearby individuals coalesce rapidly under limited dispersal and so are connected by short branch
533 lengths, while distant individuals take much longer to coalesce than they would under random
534 mating. Mutation and recombination events in our simulation both occur at a constant rate along
535 branches of the genealogy, so the genetic distance and number of recombination events separating
536 sampled individuals simply gives a noisy picture of the genealogies connecting them. Tip branches
537 (i.e., branches subtending only one individual) are then relatively short, and branches in the middle of
538 the genealogy connecting local groups of individuals relatively long, leading to the biases in the site
539 frequency spectrum shown in Figure 4.

540 The genealogical patterns introduced by limited dispersal are particularly apparent in the distribution
541 of haplotype block lengths (Figure 4). This is because identical-by-state tract lengths reflect the
542 impacts of two processes acting along the branches of the underlying genealogy – both mutation and
543 recombination – rather than just mutation as is the case when looking at the site frequency spectrum or
544 related summaries. This means that the pairwise distribution of haplotype block lengths carries with
545 it important information about genealogical variation in the population, and correlation coefficients
546 between moments of the this distribution and geographic location contain signal similar to the correlations
547 between F_{ST} or D_{xy} and geographic distance (Rousset 1997). Indeed this basic logic underlies
548 two recent studies explicitly estimating dispersal from the distribution of shared haplotype block

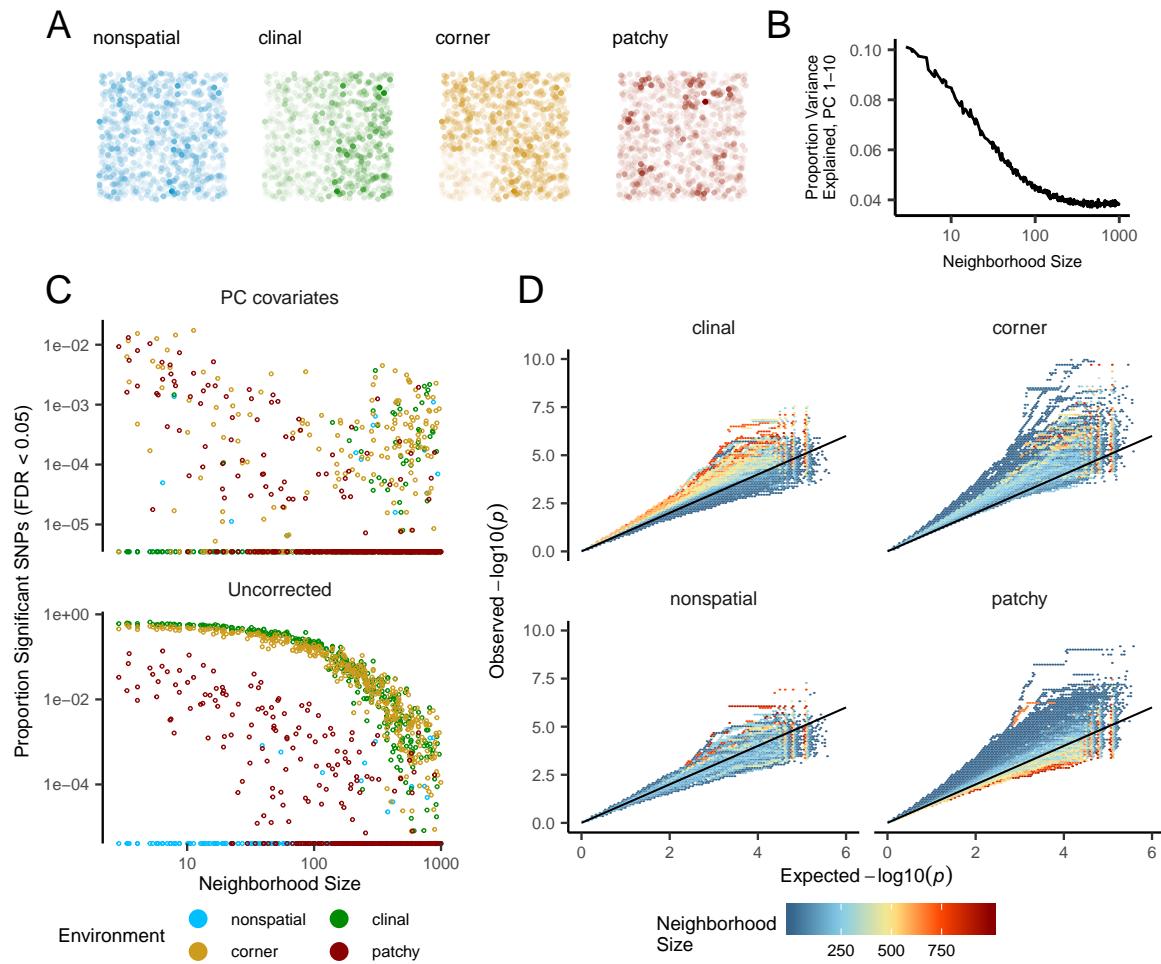


Figure 8 Impacts of spatially varying environments and isolation by distance on linear regression GWAS. Simulated quantitative phenotypes are determined only by an individual's location and the spatial distribution of environmental factors. In A we show the phenotypes and locations of sampled individuals under four environmental distributions, with transparency scaled to phenotype. As neighborhood size increases a PCA explains less of the total variation in the data (B). Spatially correlated environmental factors cause false positives at a large proportion of SNPs, which is partially but not entirely corrected by adding the first 10 PC coordinates as covariates (C). Quantile-quantile plots in D show inflation of $-\log_{10}(p)$ after PC correction across all simulations and environmental distributions, with colors scaled by the median neighborhood size in each region of q-q space.

lengths (Ringbauer *et al.* 2017; Baharian *et al.* 2016). Conversely, because haplotype-based measures of demography are particularly sensitive to variation in the underlying genealogy, inference approaches that assume random mating when analyzing the distribution of shared haplotype block lengths are likely to be strongly affected by spatial processes.

Effects of Sampling

One of the most important differences between random mating and spatial models is the effect of sampling: in a randomly mating population the spatial distribution of sampling effort has no effect on estimates of genetic variation (Table S1), but when dispersal is limited sampling strategy can compound spatial patterns in the underlying genealogy and create pervasive impacts on all downstream genetic analyses (see also Städler *et al.* (2009)). In most species, the difficulty of traveling through all parts of a species range and the inefficiency of collecting single individuals at each sampling site means that most studies follow something closest to the “point” sampling strategy we simulated, in which multiple individuals are sampled from nearby points on the landscape. For example, in ornithology a sample of 10 individuals per species per locality is a common target when collecting for natural history museums. In classical studies of *Drosophila* variation the situation is considerably worse, in which a single orchard might be extensively sampled.

When sampling is clustered at points on a landscape and dispersal is limited, the sampled individuals will be more closely related than a random set of individuals. Average coalescence times of individuals collected at a locality will then be more recent and branch lengths shorter than expected by analyses assuming random mating. This leads to fewer mutations and recombination events occurring since their last common ancestor, causing a random set of individuals to share longer average IBS tracts and have fewer nucleotide differences. For some data summaries, such as Tajima’s D , Watterson’s Θ , or the correlation coefficient between spatial distance and the count of long haplotype blocks, this can result in large differences in estimates between random and point sampling (Figure 4). Inferring underlying demographic parameters from these summary statistics – unless the nature of the sampling is somehow taken into account – will be subject to bias if sampling is not random across the landscape.

However, we observed the largest sampling effects using “midpoint” sampling. This model is meant to reflect a bias in sampling effort towards the middle of a species’ range. In empirical studies this sampling strategy could arise if, for example, researchers choose to sample the center of the range and avoid range edges to maximize probability of locating individuals during a short field season. Because midpoint sampling provides limited spatial resolution it dramatically reduces the magnitude of observed correlations between spatial and genetic distances. More surprisingly, midpoint sampling also leads to strongly positive Tajima’s D and an inflation in the proportion of heterozygous individuals in the sample – similar to the effect of sampling a single deme in an island model as reported in (Städler *et al.* 2009). This increase in observed heterozygosity appears to reflect the effects of range edges, which are a fundamental facet of spatial genetic variation. If individuals move randomly in a finite two-dimensional landscape then regions in the middle of the landscape receive migrants from all directions while those on the edge receive no migrants from at least one direction. The average number of new mutations moving into the middle of the landscape is then higher than the number moving into regions near the range edge, leading to higher heterozygosity and lower inbreeding coefficients (F_{IS}) away from range edges. Though here we used only a single parameterization of fitness decline at range edges we believe this is a general property of non-infinite landscapes as it has also been observed in previous studies simulating under lattice models (Neel *et al.* 2013; Shirk and Cushman 2014).

In summary, we recommend that empirical researchers collect individuals from across as much of the species’ range as practical, choosing samples separated by a range of spatial scales. Many summary statistics are designed for well-mixed populations, and so provide different insights into genetic variation when applied to different subsets of the population. Applied to a cluster of samples, summary statistics based on segregating sites (e.g., Watterson’s Θ and Tajima’s D), heterozygosity, or the distribution of long haplotype blocks, can be expected to depart significantly from what would be obtained from a wider distribution of samples. Comparing the results of analyses conducted on all individuals versus those limited to single individuals per locality can provide an informative contrast. Finally we wish to point out that the bias towards intermediate allele frequencies that we observe may

mean that the importance of linked selection, at least as is gleaned from the site frequency spectrum, may be systematically underestimated currently.

Demography

Previous studies have found that population structure and nonrandom sampling can create spurious signals of population bottlenecks when attempting to infer demographic history with microsatellite variation, summary statistics, or runs of homozygosity (Chikhi *et al.* 2010; Städler *et al.* 2009; Ptak and Przeworski 2002; Mazet *et al.* 2015). Here we found that methods that infer detailed population trajectories through time based on the SFS and patterns of LD across the genome are also subject to this bias, with some combinations of dispersal and sampling strategy systematically inferring deep recent population bottlenecks and overestimating ancient N_e by around a factor of 2. We were surprised to see that both stairwayplot and SMC++ can tolerate relatively strong isolation by distance – i.e., neighborhood sizes of 20 – and still perform well when averaging results across multiple simulations. Inference in populations with neighborhood sizes over 20 was relatively unbiased unless samples were concentrated in the middle of the range (Figure 7). Although median demography estimates across many independent simulations were fairly accurate, empirical work has only a single estimate to work with, and individual model fits (Figure S7) suggest that spuriously inferred population size changes and bottlenecks are common, especially at small neighborhood sizes. As we will discuss below, most empirical estimates of neighborhood size, including all estimates for human populations, are large enough that population size trajectories inferred by these approaches should not be strongly affected by spatial biases created by dispersal in continuous landscapes. In contrast, Mazet *et al.* (2015) found that varying migration rates through time could create strong biases in inferred population trajectories from an n -island model with parameters relevant for human history, suggesting that changes in migration rates through time are more likely to drive variation in inferred N_e than isolation by distance.

We found that SMC++ was more robust to the effects of space than stairwayplot, underestimating recent populations by roughly half in the worst time periods rather than nearly 10-fold as with stairwayplot. Though this degree of variation in population size is certainly meaningful in an ecological context, it is relatively minor in population genetic terms. In general methods directly assessing haplotype structure in phased data (for example, Browning and Browning (2015)) are thought to provide increased resolution for recent demographic events, but in this case the error we observed was essentially an accurate reflection of underlying genealogies in which terminal branches are anomalously short. Combined with our analysis of IBS tract length variation (Figure 5) this suggests that haplotype-based methods are likely to be affected by similar biases.

A more worrying pattern was the high level of variance in inferred N_e trajectories for individual model fits using these methods, which was highest in simulations with the smallest neighborhood size (Figure 7, Figure S7). This suggests that, at a minimum, researchers working with empirical data should replicate analyses multiple times and take a rolling average if model fits are inconsistent across runs. Splitting samples and running replicates on separate subsets – the closest an empirical study can come to our design of averaging the results from multiple simulations – may also alleviate this issue.

Our analysis suggests that many empirical analyses of population size history using methods like SMC++ are robust to error caused by spatial structure within continuous landscapes. Inferences drawn from static SFS-based methods like stairwayplot should be treated with caution when there are signs of isolation by distance in the underlying data (for example, if a regression of F_{ST} against the logarithm of geographic distance has a significantly positive slope), and in particular an inference of population bottlenecks in the last 1000 years should be discounted if sampling is clustered, but estimates of deeper time patterns are likely to be fairly accurate. The biases in the SFS and haplotype structure identified above (see also Wakeley 1999; Chikhi *et al.* 2010; Städler *et al.* 2009) are apparently small enough that they fall within the range of variability regularly inferred by these approaches, at least on datasets of the size we simulated.

GWAS

Spatial structure is particularly challenging for genome-wide association studies, because the effects of dispersal on genetic variation are compounded by spatial variation in the environment (Mathieson

and McVean 2012). Spatially restricted mate choice and dispersal causes variation in allele frequencies across the range of a species. If environmental factors affecting the phenotype of interest also vary over space, then allele frequencies and environmental exposures will covary over space. In this scenario an uncorrected GWAS will infer genetic associations with a purely environmental phenotype at any site in the genome that is differentiated over space, and the relative degree of bias will be a function of the degree of covariation in allele frequencies and the environment (i.e., Figure 8C, bottom panel). This pattern has been demonstrated in a variety of simulation and empirical contexts (Price *et al.* 2006; Yu *et al.* 2005; Young *et al.* 2018; Mathieson and McVean 2012; Kang *et al.* 2008, 2010; Bulik-Sullivan *et al.* 2015; Berg *et al.* 2018; Sohail *et al.* 2018).

Incorporating PC positions as covariates in a linear-regression GWAS (Price *et al.* 2006) is designed to address this challenge by regressing out a baseline level of “average” differentiation. In essence, a PC-corrected GWAS asks “what regions of the genome are more associated with this phenotype than the average genome-wide association observed across populations?” In our simulations, we observed that this procedure can fail under a variety of circumstances. If dispersal is limited and environmental variation is clustered in space (i.e., corner or patchy distributions in our simulations), PCA positions fail to capture the fine-scale spatial structure required to remove all signals of association. Conversely, as dispersal increases, PCA loses power to describe population structure before spatial mixing breaks down the relationship between genotype and the environment. These effects were observed with all spatially correlated environmental patterns, but were particularly pronounced if environmental effects are concentrated in one region, as was also found by Mathieson and McVean (2012). Though increasing the number of PC axes used in the analysis may reduce the false-positive rate, this may also decrease the power of the test to detect truly causal alleles (Lawson *et al.* 2019).

In this work we simulated a single chromosome with size roughly comparable to one human chromosome. If we scale the number of false-positive associations identified in our analyses to a GWAS conducted on whole-genome data from humans, we would expect to see several thousand weak false-positive associations after PC corrections in a population with neighborhood sizes up to at least 1000 (which should include values appropriate for many human populations). Notably, very few of the spurious associations we identified would be significant at a conservative Bonferroni-adjusted *p*-value cutoff (see Figure S8). This suggests that GWAS focused on finding strongly associated alleles for traits controlled by a limited number of variants in the genome are likely robust to the impacts of continuous spatial structure. However, methods that analyze the combined effects of thousands or millions of weakly associated variants such as polygenic risk scores (Khera *et al.* 2018) are likely to be affected by subtle population structure. Indeed as recently identified in studies of genotype associations for human height in Europe (Berg *et al.* 2018; Sohail *et al.* 2018), PC regression GWAS in modern human populations do include residual signal of population structure in large-scale analyses of polygenic traits. When attempting to make predictions across populations with different environmental exposures, polygenic risk scores affected by population structure can be expected to offer low predictive power, as was shown in a recent study finding lower performance outside European populations (Martin *et al.* 2019).

In summary, spatial covariation in population structure and the environment confounds the interpretation of GWAS *p*-values, and correction using principal components is insufficient to fully separate these signals for polygenic traits under a variety of environmental and population parameter regimes. Other GWAS methods such as mixed models (Kang *et al.* 2008) may be less sensitive to this confounding, but there is no obvious reason that this should be so. One approach to estimating the degree of bias in GWAS caused by population structure is LD score regression (Bulik-Sullivan *et al.* 2015). Though this approach appears to work well in practice, its interpretation is not always straightforward and it is likely biased by the presence of linked selection (Berg *et al.* 2018). In addition, we observed that in many cases the false-positive SNPs we identified appeared to be concentrated in LD peaks similar to those expected from truly causal sites (Figure S8), which may confound LD score regression.

We suggest a straightforward alternative for species in which the primary axes of population differentiation is space (note this is likely not the case for some modern human populations): run a GWAS with spatial coordinates as phenotypes and check for *p*-value inflation or significant associations.

705 If significant associations with sample locality are observed after correcting for population structure,
 706 the method is sensitive to false positives induced by spatial structure. This is essentially the approach
 707 taken in our “clinal” model (though we add normally distributed noise to our phenotypes). This
 708 approach has recently been taken with polygenic scores for UK Biobank samples in Haworth *et al.*
 709 (2019), finding that scores are correlated with birth location even in this relatively homogenous sample
 710 . Of course, it is possible that genotypes indirectly affect individual locations by adjusting organismal
 711 fitness and thus habitat selection across spatially varying environments, but we believe that this
 712 hypothesis should be tested against a null of stratification bias inflation rather than accepted as true
 713 based on GWAS results.

Table 1 Neighborhood size estimates from empirical studies.

Species	Description	Neighborhood Size	Method	Citation
<i>Ipomopsis aggregata</i>	flowering plant	12.60 - 37.80	Genetic	(Campbell and Dooley 1992)
<i>Borreria frutescens</i>	salt marsh plant	20 - 30	Genetic+Survey	(Antlfinger 1982)
<i>Oreamnos americanus</i>	mountain goat	36 - 100	Genetic	(Shirk and Cushman 2014)
<i>Homo sapiens</i>	Gainj- and Kalam-speaking people, Papua New Guinea	40 - 213	Genetic	(Rousset 1997)
<i>Formica sp.</i>	colonial ants	50 - 100	Genetic	(Pamilo 1983)
<i>Astrocaryum mexicanum</i>	palm tree	102 - 895	Genetic+survey	(Eguiarte <i>et al.</i> 1993)
<i>Spermophilus mollis</i>	ground squirrel	204 - 480	Genetic+Survey	(Antolin <i>et al.</i> 2001)
<i>Sceloporus olivaceus</i>	lizard	225 - 270	Survey	(Kerster 1964)
<i>Dieffenbachia longispatha</i>	beetle-pollinated colonial herb	227 - 611	Survey	(Young 1988)
<i>Aedes aegypti</i>	Yellow-fever mosquito	268	Genetic	(Jasper <i>et al.</i> 2019)
<i>Homo sapiens</i>	Gainj- and Kalam-speaking people, Papua New Guinea	410	Survey	(Rousset 1997)
<i>Quercus laevis</i>	Oak tree	> 440	Genetic	(Berg and Hamrick 1995)
<i>Drosophila pseudoobscura</i>	fruit fly	500 - 1,000	Survey+Crosses	(Wright 1946)
<i>Homo sapiens</i>	POPRES data NE Europe	1,342 - 5,425	Genetic	(Ringbauer <i>et al.</i> 2017)
<i>Bebicium vittatum</i>	intertidal snail	240,000	Survey	(Rousset 1997)
<i>Bebicium vittatum</i>	intertidal snail	360,000	Genetic	(Rousset 1997)

714 **Where are natural populations on this spectrum?**

715 For how much of the tree of life do spatial patterns circumscribe genomic variation? In Table 1 we
 716 gathered estimates of neighborhood size from a range of organisms to get an idea of how likely
 717 dispersal is to play an important role in patterns of variation. These values should be compared to
 718 our simulation results with some caution, as though we expect neighborhood size to have similar
 719 effects across species of varying global N_e (Wright 1946), in our study we evaluated only a relatively

720 small population of $\approx 10,000$. In addition, these empirical examples are likely biased towards small-
721 neighborhood species (because few studies have quantified neighborhood size in species with very
722 high dispersal or population density). However, from the available data we find that neighborhood
723 sizes in the range we simulated are fairly common across a range of taxa. At the extreme low end of
724 empirical neighborhood size estimates we see some flowering plants, large mammals, and colonial
725 insects like ants. Species such as this have neighborhood size estimates small enough that spatial
726 processes are likely to strongly influence inference. These include some human populations such as
727 the Gainj- and Kalam-speaking people of Papua New Guinea, in which the estimated neighborhood
728 sizes in (Rousset 1997) range from 40 to 410 depending on the method of estimation. Many more
729 species occur in a middle range of neighborhood sizes between 100 and 1000 – a range in which spatial
730 processes play a minor role in our analyses under random spatial sampling but are important when
731 sampling of individuals in space is clustered. Surprisingly, even some flying insects with huge census
732 population sizes fall in this group, including fruit flies (*D. melanogaster*) and mosquitoes (*A. aegypti*).
733 Last, many species likely have neighborhood sizes much larger than we simulated, including modern
734 humans in northeastern Europe (Ringbauer *et al.* 2017). For these species demographic inference
735 and summary statistics are likely to reflect minimal bias from spatial effects as long as dispersal is
736 truly continuous across the landscape. While that is so we caution that association studies in which
737 the effects of population structure are confounded with spatial variation in the environment are still
738 sensitive to dispersal even at these large neighborhood sizes.

739 **Future Directions and Limitations**

740 As we have shown, a large number of population genetic summary statistics contain information about
741 spatial population processes. We imagine that combinations of such summaries might be sufficient
742 for the construction of supervised machine learning regressors (e.g., Schrider and Kern 2018) for the
743 accurate estimation of dispersal from genetic data. Indeed, Ashander *et al.* (2018) found that inverse
744 interpolation on a vector of summary statistics provided a powerful method of estimating dispersal
745 distances. Expanding this approach to include the haplotype-based summary statistics studied here
746 and applying machine learning regressors built for general inference of nonlinear relationships from
747 high-dimensional data may allow precise estimation of spatial parameters under a range of complex
748 models.

749 One facet of spatial variation that we did not address in this study is the confounding of dispersal
750 and population density implicit in the definition of Wright's neighborhood size. Our simulations were
751 run under constant densities, but Guindon *et al.* (2016) and Ringbauer *et al.* (2017) have shown that
752 these parameters are identifiable under some continuous models. Similarly, though the scaling effects
753 of dispersal we show in Figure 4 should occur in populations of any total size, other aspects such as
754 the number of segregating sites are also likely affected by the total landscape size (and so total census
755 N). Much additional work remains to be done to better understand how these parameters interact to
756 shape genetic variation in continuous space, which we leave to future studies.

757 Though our simulation allows incorporation of realistic demographic and spatial processes, it
758 is inevitably limited by the computational burden of tracking tens or hundreds of thousands of
759 individuals in every generation. In particular, computations required for mate selection and spatial
760 competition scale approximately with the product of the total census size and the neighborhood
761 size and so increase rapidly for large populations and dispersal distances. The reverse-time model
762 of continuous space evolution described by Barton *et al.* (2010) and implemented by Kelleher *et al.*
763 (2014) allows exploration of parameter regimes with population and landscape sizes more directly
764 comparable to empirical cases like humans. Alternatively, implementation of parallelized calculations
765 may allow progress with forward-time simulations.

766 Finally, we believe that the difficulties in correcting for population structure in continuous popula-
767 tions using principal components analysis or similar decompositions is a difficult issue, well worth
768 considering on its own. How can we best avoid spurious correlations while correlating genetic and
769 phenotypic variation without underpowering the methods? Perhaps optimistically, we posit that
770 process-driven descriptions of ancestry and/or more generalized unsupervised methods may be able
771 to better account for carry out this task.

772 **Data Availability**

773 Scripts used for all analyses and figures are available at <https://github.com/petrelharp/spaceness>.

774 **Acknowledgements**

775 We thank Brandon Cooper, Matt Hahn, Doc Edge, and others for reading and thinking about this
776 manuscript. CJB and ADK were supported by NIH award R01GM117241.

777 **Literature Cited**

- 778 Aguillon, S. M., J. W. Fitzpatrick, R. Bowman, S. J. Schoech, A. G. Clark, *et al.*, 2017 Deconstructing
779 isolation-by-distance: The genomic consequences of limited dispersal. *PLOS Genetics* **13**: 1–27.
- 780 Al-Asadi, H., D. Petkova, M. Stephens, and J. Novembre, 2019 Estimating recent migration and
781 population-size surfaces. *PLoS genetics* **15**: e1007908.
- 782 Allee, W. C., O. Park, A. E. Emerson, T. Park, K. P. Schmidt, *et al.*, 1949 Principles of animal ecology.
783 Technical report, Saunders Company Philadelphia, Pennsylvania, USA.
- 784 Antlfinger, A. E., 1982 Genetic neighborhood structure of the salt marsh composite, *Borrichia frutescens*.
785 *Journal of Heredity* **73**: 128–132.
- 786 Antolin, M. F., B. V. Horne, M. D. Berger, Jr., A. K. Holloway, J. L. Roach, *et al.*, 2001 Effective population
787 size and genetic structure of a piute ground squirrel (*Spermophilus mollis*) population. *Canadian
788 Journal of Zoology* **79**: 26–34.
- 789 Antonovics, J. and D. A. Levin, 1980 The ecological and genetic consequences of density-dependent
790 regulation in plants. *Annual Review of Ecology and Systematics* **11**: 411–452.
- 791 Ashander, J., P. Ralph, E. McCartney-Melstad, and H. B. Shaffer, 2018 Demographic inference in a
792 spatially-explicit ecological model from genomic data: a proof of concept for the mojave desert
793 tortoise. *bioRxiv* .
- 794 Baharian, S., M. Barakatt, C. R. Gignoux, S. Shringarpure, J. Errington, *et al.*, 2016 The great migration
795 and african-american genomic diversity. *PLOS Genetics* **12**: 1–27.
- 796 Barton, N. H., F. Depaulis, and A. M. Etheridge, 2002 Neutral evolution in spatially continuous
797 populations. *Theoretical Population Biology* **61**: 31–48.
- 798 Barton, N. H., J. Kelleher, and A. M. Etheridge, 2010 A new model for extinction and recolonization in
799 two dimensions: Quantifying phylogeography. *Evolution* **64**: 2701–2715.
- 800 Benjamini, Y. and D. Yekutieli, 2001 The control of the false discovery rate in multiple testing under
801 dependency. *The Annals of Statistics* **29**: 1165–1188.
- 802 Berg, E. E. and J. L. Hamrick, 1995 Fine-scale genetic structure of a turkey oak forest. *Evolution* **49**:
803 110–120.
- 804 Berg, J. J., A. Harpak, N. Sinnott-Armstrong, A. M. Joergensen, H. Mostafavi, *et al.*, 2018 Reduced
805 signal for polygenic adaptation of height in uk biobank. *bioRxiv* .
- 806 Bolker, B. and S. W. Pacala, 1997 Using moment equations to understand stochastically driven spatial
807 pattern formation in ecological systems. *Theoretical Population Biology* **52**: 179 – 197.
- 808 Bolker, B. M., S. W. Pacala, and C. Neuhauser, 2003 Spatial dynamics in model plant communities:
809 What do we really know? *The American Naturalist* **162**: 135–148, PMID: 12858259.
- 810 Browning, S. R. and B. L. Browning, 2015 Accurate non-parametric estimation of recent effective
811 population size from segments of identity by descent. *The American Journal of Human Genetics* **97**:
812 404–418.
- 813 Bulik-Sullivan, B. K., P.-R. Loh, H. K. Finucane, S. Ripke, J. Yang, *et al.*, 2015 Ld score regression
814 distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics* **47**: 291 EP –.
- 815 Campbell, D. R. and J. L. Dooley, 1992 The spatial scale of genetic differentiation in a hummingbird-
816 pollinated plant: Comparison with models of isolation by distance. *The American Naturalist* **139**:
817 735–748.
- 818 Champer, J., I. Kim, S. E. Champer, A. G. Clark, and P. W. Messer, 2019 Suppression gene drive in
819 continuous space can result in unstable persistence of both drive and wild-type alleles. *bioRxiv* .

- 821 Chapman, N. H. and E. A. Thompson, 2002 The effect of population history on the lengths of ancestral
 822 chromosome segments. *Genetics* **162**: 449–458.
- 823 Chikhi, L., V. C. Sousa, P. Luisi, B. Goossens, and M. A. Beaumont, 2010 The confounding effects of
 824 population structure, genetic diversity and the sampling scheme on the detection and quantification
 825 of population size changes. *Genetics* **186**: 983–995.
- 826 Crawley, M. J., 1990 The population dynamics of plants. *Philosophical Transactions of the Royal Society*
 827 of London. Series B: Biological Sciences **330**: 125–140.
- 828 Durrett, R. and S. Levin, 1994 The importance of being discrete (and spatial). *Theoretical Population*
 829 *Biology* **46**: 363–394.
- 830 Eguiarte, L. E., A. Búrquez, J. Rodríguez, M. Martínez-Ramos, J. Sarukhán, *et al.*, 1993 Direct and
 831 indirect estimates of neighborhood and effective population size in a tropical palm, *astrocaryum*
 832 *mexicanum*. *Evolution* **47**: 75–87.
- 833 Epperson, B., 2003 *Geographical Genetics*. Monographs in Population Biology, Princeton University
 834 Press.
- 835 Felsenstein, J., 1975 A pain in the torus: Some difficulties with models of isolation by distance. *The*
 836 *American Naturalist* **109**: 359–368.
- 837 Fournier, N. and S. Méléard, 2004 A microscopic probabilistic description of a locally regulated
 838 population and macroscopic approximations. *The Annals of Applied Probability* **14**: 1880–1919.
- 839 Fox, J. and S. Weisberg, 2011 *An R Companion to Applied Regression*. Sage, Thousand Oaks CA, second
 840 edition.
- 841 Garcia, J. and C. Quintana-Domeque, 2006 The evolution of adult height in europe: A brief note.
 842 Working Paper .
- 843 Garcia, J. and C. Quintana-Domeque, 2007 The evolution of adult height in europe: A brief note.
 844 *Economics & Human Biology* **5**: 340 – 349.
- 845 Garud, N. R., P. W. Messer, E. O. Buzbas, and D. A. Petrov, 2015 Recent selective sweeps in north
 846 american *drosophila melanogaster* show signatures of soft sweeps. *PLOS Genetics* **11**: 1–32.
- 847 Griffiths, R., S. Tavaré, *et al.*, 1999 The ages of mutations in gene trees. *The Annals of Applied Probability*
 848 **9**: 567–590.
- 849 Guindon, S., H. Guo, and D. Welch, 2016 Demographic inference under the coalescent in a spatial
 850 continuum. *Theoretical population biology* **111**: 43–50.
- 851 Haller, B. C., J. Galloway, J. Kelleher, P. W. Messer, and P. L. Ralph, 2019 Tree-sequence recording
 852 in SLiM opens new horizons for forward-time simulation of whole genomes. *Molecular Ecology*
 853 *Resources* **19**: 552–566.
- 854 Haller, B. C. and P. W. Messer, 2019 Slim 3: Forward genetic simulations beyond the wright-fisher
 855 model. *Molecular biology and evolution* **36**: 632–637.
- 856 Harris, K. and R. Nielsen, 2013 Inferring demographic history from a spectrum of shared haplotype
 857 lengths. *PLOS Genetics* **9**: 1–20.
- 858 Haworth, S., R. Mitchell, L. Corbin, K. H. Wade, T. Dudding, *et al.*, 2019 Apparent latent structure
 859 within the uk biobank sample has implications for epidemiological analysis. *Nature communications*
 860 **10**: 333.
- 861 Huillet, T. and M. Möhle, 2011 On the extended Moran model and its relation to coalescents with
 862 multiple collisions. *Theoretical Population Biology* pp. –.
- 863 Jasper, M., T. Schmidt, N. Ahmad, S. Sinkins, and A. Hoffmann, 2019 A genomic approach to inferring
 864 kinship reveals limited intergenerational dispersal in the yellow fever mosquito. *bioRxiv* .
- 865 Jay, F., P. Sjödin, M. Jakobsson, and M. G. Blum, 2012 Anisotropic Isolation by Distance: The Main
 866 Orientations of Human Genetic Differentiation. *Molecular Biology and Evolution* **30**: 513–525.
- 867 Kang, H. M., J. H. Sul, S. K. Service, N. A. Zaitlen, S.-y. Kong, *et al.*, 2010 Variance component model to
 868 account for sample structure in genome-wide association studies. *Nature Genetics* **42**: 348 EP –.
- 869 Kang, H. M., N. A. Zaitlen, C. M. Wade, A. Kirby, D. Heckerman, *et al.*, 2008 Efficient control of
 870 population structure in model organism association mapping. *Genetics* **178**: 1709–1723.
- 871 Kelleher, J., A. Etheridge, and N. Barton, 2014 Coalescent simulation in continuous space: Algorithms
 872 for large neighbourhood size. *Theoretical Population Biology* **95**: 13 – 23.
- 873 Kelleher, J., A. M. Etheridge, and G. McVean, 2016 Efficient coalescent simulation and genealogical

- analysis for large sample sizes. *PLoS Comput Biol* **12**: 1–22.
- Kelleher, J., K. R. Thornton, J. Ashander, and P. L. Ralph, 2018 Efficient pedigree recording for fast population genetics simulation. *PLOS Computational Biology* **14**: 1–21.
- Kerster, H. W., 1964 Neighborhood size in the rusty lizard, *sceloporus olivaceus*. *Evolution* **18**: 445–457.
- Khera, A. V., M. Chaffin, K. G. Aragam, M. E. Haas, C. Roselli, *et al.*, 2018 Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nature Genetics* **50**: 1219–1224.
- Kingman, J., 1982 The coalescent. *Stochastic Processes and their Applications* **13**: 235 – 248.
- Law, R., D. J. Murrell, and U. Dieckmann, 2003 Population growth in space and time: Spatial logistic equations. *Ecology* **84**: 252–262.
- Lawson, D. J., N. M. Davies, S. Haworth, B. Ashraf, L. Howe, *et al.*, 2019 Is population structure in the genetic biobank era irrelevant, a challenge, or an opportunity? *Human Genetics* .
- Liu, X. and Y.-X. Fu, 2015 Exploring population size changes using snp frequency spectra. *Nature Genetics* **47**: 555 EP –.
- Lloyd, M., 1967 'Mean crowding'. *Journal of Animal Ecology* **36**: 1–30.
- Lundgren, E. and P. L. Ralph, 2018 Are populations like a circuit? The relationship between isolation by distance and isolation by resistance. *bioRxiv* .
- Martin, A. R., M. Kanai, Y. Kamatani, Y. Okada, B. M. Neale, *et al.*, 2019 Clinical use of current polygenic risk scores may exacerbate health disparities. *Nature Genetics* **51**: 584–591.
- Maruyama, T., 1972 Rate of decrease of genetic variability in a two-dimensional continuous population of finite size. *Genetics* **70**: 639–651.
- Mathieson, I. and G. McVean, 2012 Differential confounding of rare and common variants in spatially structured populations. *Nature Genetics* **44**: 243 EP –.
- Mazet, O., W. Rodríguez, S. Grusea, S. Boitard, and L. Chikhi, 2015 On the importance of being structured: instantaneous coalescence rates and human evolution—lessons for ancestral population size inference? *Heredity* **116**: 362 EP –.
- Miles, A. and N. Harding, 2017 *cghg/scikit-allel*: v1.1.8.
- Neel, M. C., K. McKelvey, N. Ryman, M. W. Lloyd, R. Short Bull, *et al.*, 2013 Estimation of effective population size in continuously distributed populations: there goes the neighborhood. *Heredity* **111**: 189 EP –.
- Novembre, J. and M. Slatkin, 2009 Likelihood-based inference in isolation-by-distance models using the spatial distribution of low-frequency alleles. *Evolution* **63**: 2914–2925.
- Pamilo, P., 1983 Genetic differentiation within subdivided populations of *formica* ants. *Evolution* **37**: 1010–1022.
- Patterson, N., A. L. Price, and D. Reich, 2006 Population structure and eigenanalysis. *PLOS Genetics* **2**: 1–20.
- Price, A. L., N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, *et al.*, 2006 Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics* **38**: 904 EP –.
- Pritchard, J. K., M. Stephens, and P. Donnelly, 2000 Inference of population structure using multilocus genotype data. *Genetics* **155**: 945–959.
- Ptak, S. E. and M. Przeworski, 2002 Evidence for population growth in humans is confounded by fine-scale population structure. *Trends in Genetics* **18**: 559–563.
- Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. A. Ferreira, *et al.*, 2007 Plink: A tool set for whole-genome association and population-based linkage analyses. *The American Journal of Human Genetics* **81**: 559 – 575.
- R Core Team, 2018 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Ralph, P. and G. Coop, 2013 The geography of recent genetic ancestry across Europe. *PLoS Biol* **11**: e1001555.
- Ringbauer, H., G. Coop, and N. H. Barton, 2017 Inferring recent demography from isolation by distance of long shared sequence blocks. *Genetics* **205**: 1335–1351.
- Robledo-Arnuncio, J. J. and F. Rousset, 2010 Isolation by distance in a continuous population under

- 927 stochastic demographic fluctuations. *Journal of Evolutionary Biology* **23**: 53–71.
- 928 Rossine, F. W. S., 2014 *Espaço e diversificação: uma perspectiva teórica*. Master's dissertation in ecologia:
929 Ecossistemas terrestres e aquáticos, University of São Paulo, São Paulo : Instituto de Biociências.
- 930 Rousset, F., 1997 Genetic differentiation and estimation of gene flow from F-statistics under isolation
931 by distance. *Genetics* **145**: 1219–1228.
- 932 Rousset, F. and R. Leblois, 2011 Likelihood-based inferences under isolation by distance: Two-
933 dimensional habitats and confidence intervals. *Molecular Biology and Evolution* **29**: 957–973.
- 934 Sawyer, S., 1977 On the past history of an allele now known to have frequency p. *Journal of Applied
935 Probability* **14**: 439–450.
- 936 Schiffels, S. and R. Durbin, 2014 Inferring human population size and separation history from multiple
937 genome sequences. *Nature Genetics* **46**: 919 EP –.
- 938 Schrider, D. R. and A. D. Kern, 2018 Supervised machine learning for population genetics: A new
939 paradigm. *Trends in Genetics* **34**: 301 – 312.
- 940 Sharbel, T. F., B. Haubold, and T. Mitchell-Olds, 2000 Genetic isolation by distance in arabidopsis
941 thaliana: biogeography and postglacial colonization of europe. *Molecular Ecology* **9**: 2109–2118.
- 942 Sheehan, S., K. Harris, and Y. S. Song, 2013 Estimating variable effective population sizes from multiple
943 genomes: A sequentially markov conditional sampling distribution approach. *Genetics* **194**: 647–662.
- 944 Shirk, A. J. and S. A. Cushman, 2014 Spatially-explicit estimation of wright's neighborhood size in
945 continuous populations. *Frontiers in Ecology and Evolution* **2**: 62.
- 946 Slatkin, M. and N. H. Barton, 1989 A comparison of three indirect methods for estimating average
947 levels of gene flow. *Evolution* **43**: 1349–1368.
- 948 Sohail, M., R. M. Maier, A. Ganna, A. Bloemendal, A. R. Martin, *et al.*, 2018 Signals of polygenic
949 adaptation on height have been overestimated due to uncorrected population structure in genome-
950 wide association studies. *bioRxiv* .
- 951 St. Onge, K. R., A. E. Palmé, S. I. Wright, and M. Lascoux, 2012 Impact of sampling schemes on
952 demographic inference: An empirical study in two species with different mating systems and
953 demographic histories. *G3: Genes, Genomes, Genetics* **2**: 803–814.
- 954 Städler, T., B. Haubold, C. Merino, W. Stephan, and P. Pfaffelhuber, 2009 The impact of sampling
955 schemes on the site frequency spectrum in nonequilibrium subdivided populations. *Genetics* **182**:
956 205–216.
- 957 Tajima, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism.
958 *Genetics* **123**: 585–595.
- 959 Terhorst, J., J. A. Kamm, and Y. S. Song, 2016 Robust and scalable inference of population history from
960 hundreds of unphased whole genomes. *Nature Genetics* **49**: 303 EP –.
- 961 Turchin, M. C., C. W. Chiang, C. D. Palmer, S. Sankararaman, D. Reich, *et al.*, 2012 Evidence of
962 widespread selection on standing variation in europe at height-associated snps. *Nature Genetics* **44**:
963 1015 EP –.
- 964 Wahlund, S., 1928 Zusammensetzung von populationen und korrelationserscheinungen vom stand-
965 punkt der vererbungslehre aus betrachtet. *Hereditas* **11**: 65–106.
- 966 Wakeley, J., 1999 Nonequilibrium migration in human history. *Genetics* **153**: 1863–1871.
- 967 Wakeley, J., 2009 *Coevalent Theory, an Introduction*. Roberts and Company, Greenwood Village, CO.
- 968 Wakeley, J. and T. Takahashi, 2003 Gene genealogies when the sample size exceeds the effective size of
969 the population. *Mol Biol Evol* **20**: 208–213.
- 970 Wickham, H., 2016 *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- 971 Wilke, C. O., 2019 *cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'*. R package version
972 0.9.4.
- 973 Wilkins, J. F., 2004a A separation-of-timescales approach to the coalescent in a continuous population.
974 *Genetics* **168**: 2227–2244.
- 975 Wilkins, J. F., 2004b A separation-of-timescales approach to the coalescent in a continuous population.
976 *Genetics* **168**: 2227–2244.
- 977 Wilkins, J. F. and J. Wakeley, 2002 The coalescent in a continuous, finite, linear population. *Genetics*
978 **161**: 873–888.
- 979 Wright, S., 1931 Evolution in mendelian populations. *Genetics* **16**: 97.

- 980 Wright, S., 1943 Isolation by distance. *Genetics* **28**: 114–138.
- 981 Wright, S., 1946 Isolation by distance under diverse systems of mating. *Genetics* **31**: 336.
- 982 Young, A. I., M. L. Frigge, D. F. Gudbjartsson, G. Thorleifsson, G. Bjornsdottir, *et al.*, 2018 Relatedness
983 disequilibrium regression estimates heritability without environmental bias. *Nature Genetics* **50**:
984 1304–1310.
- 985 Young, H. J., 1988 Neighborhood size in a beetle pollinated tropical aroid: effects of low density and
986 asynchronous flowering. *Oecologia* **76**: 461–466.
- 987 Yu, J., G. Pressoir, W. H. Briggs, I. Vroh Bi, M. Yamasaki, *et al.*, 2005 A unified mixed-model method for
988 association mapping that accounts for multiple levels of relatedness. *Nature Genetics* **38**: 203 EP –.
- 989 Zähle, I., J. T. Cox, and R. Durrett, 2005 The stepping stone model. II. Genealogies and the infinite sites
990 model. *Ann. Appl. Probab.* **15**: 671–699.

991 Appendix 1

992

993 Comparisons with Stepping-Stone Models

994 We checked that our model produces reasonable results by comparing it to a reverse-time stepping-
995 stone model implemented in msprime (Kelleher *et al.* 2016). In this class of models we imagine an $n \times n$
996 grid of populations exchanging migrants with neighboring populations at rate m . If these models are
997 good approximations of the continuous case we expect that results will converge as $n \rightarrow \infty$, so we
998 ran simulations while varying n from 5 to 50 (Table A1). To compare with continuous models we first
999 distributed the same "effective" number of individuals across the landscape as in our continuous-
1000 space simulations (≈ 6100 , estimated from θ_π of random-mating continuous-space simulations). We
1001 then approximate the mean per-generation dispersal distance σ given a total landscape width W as
1002 the product of the probability of an individual being a migrant and the distance traveled by migrants:
1003 $\sigma \approx 4m(W/n)$. We ran 500 simulations per n while sampling σ from $U(0.2, 4)$. We then randomly
1004 selected 60 diploid individuals from each simulation (approximating diploidy by combining pairs of
1005 chromosomes with contiguous indices within demes) and calculated a set of six summary statistics
1006 using the scripts described in the summary statistics portion of the main text.

demes per side (n)	N_e per deme	samples per deme
5	244	20
10	61	10
20	15.25	2
50	2.44	1

Table A1 stepping-stone simulation parameters

1007 In general we find many of the qualitative trends are similar among continuous and stepping-stone
1008 models and that, in most cases, statistics from stepping-stone models approach the continuous model
1009 as the resolution of the grid increases. For example, θ_π is inflated at low neighborhood sizes (i.e. low
1010 m), and the extent of the inflation increases to approach the continuous case as the resolution of the
1011 landscape increases. Similar patterns are observed for F_{is} and observed heterozygosity. However, θ_W
1012 behaves differently, with increased grid resolution leading to lower values. This in turn drives an even
1013 more positive Tajima's D in grid simulations at small neighborhood sizes.

1014 These differences relative to our continuous model mainly reflect two shortcomings of the reverse-
1015 time stepping stone model. If we simulate a coarse grid with relatively large populations in each
1016 deme, we cannot accurately capture the dynamics of small neighborhood sizes because mating within
1017 each deme remains random regardless of the migration rate connecting demes. This likely explains

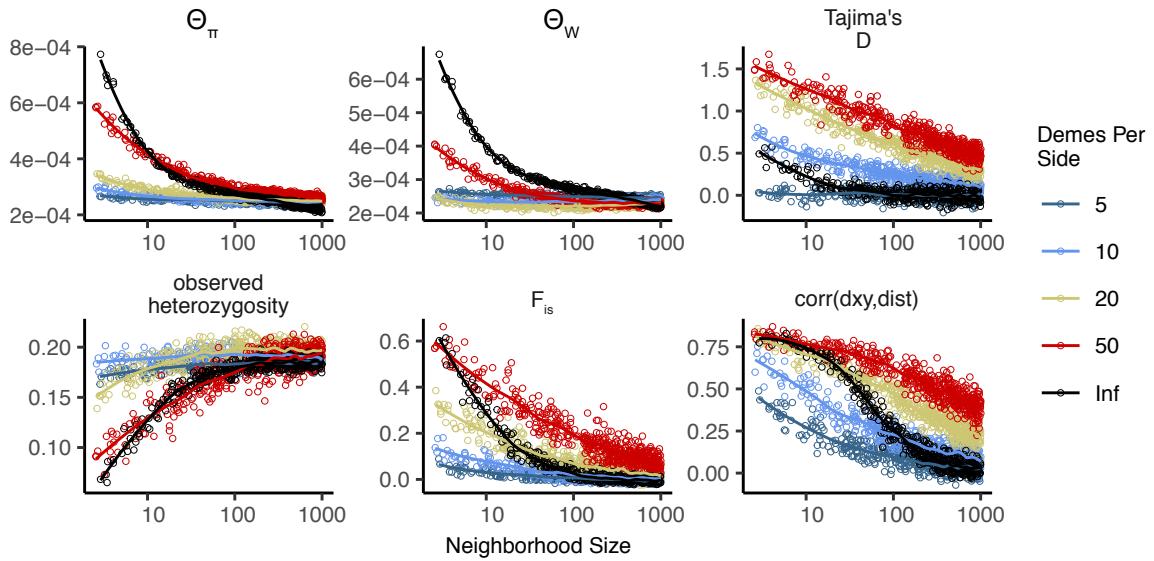


Figure A1 Summary statistics for 2-dimensional coalescent stepping-stone models with fixed total N_e and varying numbers of demes per side. The black "infinite" points are from our forward-time continuous space model. Inter-deme migration rates are related to σ as described above.

the trends in π , observed heterozygosity, and F_{is} . However increasing the number of demes while holding the total number of individuals constant results in small within-deme populations for which even the minimum sample size of 1 approaches the local N_e (Table A1). This results in an excess of short terminal branches in the coalescent tree, which decreases the total branch length and leads to fewer segregating sites, deflated θ_W , and inflated Tajima's D . Overall then our continuous model reproduces important features of spatial structure approximated by reverse-time stepping-stone models at moderate neighborhood sizes while avoiding some artifacts caused binning the landscape into discrete demes.

Demographic model

We chose our demographic model so that every individual has on average $1/L$ offspring each time step, and if the local population density of an individual is n , then their probability of survival until the next time step is (equation (1)):

$$p = \min \left(0.95, \frac{1}{1 + n/(K(1+L))} \right). \quad (3)$$

We capped survival at 0.95 so that we would not have exceptionally long-lived individuals in sparsely populated areas – otherwise, an isolated individual might live for a very long time. Since $1 - p \approx n/(K(1+L))$, mortality goes up roughly linearly with the number of neighbors (on a scale given by K), as would be obtained if, for instance, mortality is due to agonistic interactions. Ignoring in/outmigration, a region is at demographic equilibrium if the per-capita probability of death is equal to the birth rate, i.e., if $1 - p = 1/L$. (Note that there is no effect of age in the model, which would make the analysis more complicated.) Solving this for n , we get that in a well-mixed population, the equilibrium density should be around

$$n = K \frac{L+1}{L-1} \quad (4)$$

1039 individuals per unit area. At this density, the per-capita death rate is $1/L$, so the mean lifetime is L .
1040 This equilibrium density is *not* K , but (since $L = 4$) is two-thirds larger. However, in practice this model
1041 leads to a total population size which is around K multiplied by total geographic area (but which
1042 depends on σ , as discussed above). The main reason for this is that since offspring tend to be near
1043 their parents, individuals tend to be “clumped”, and so experience a higher average density than the
1044 “density” one would compute by dividing census size by geographic area (Lloyd 1967). To maintain a
1045 constant expected total population size would require making (say) K depend on σ ; however, typical
1046 local population densities might then be more dissimilar.

¹⁰⁴⁷ **Supplementary Figures and Tables**

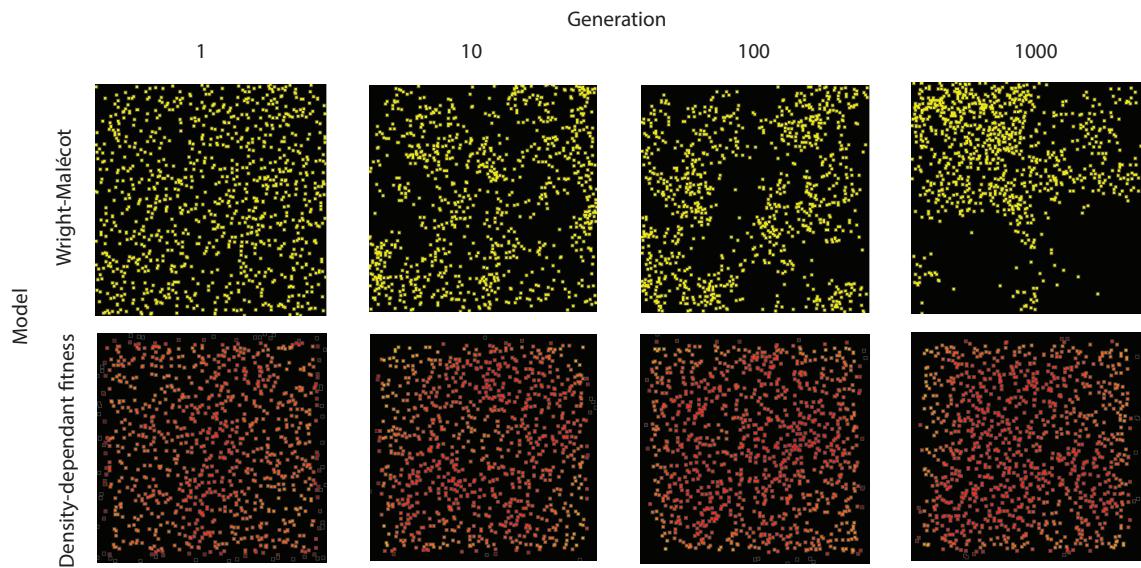


Figure S1 Maps of individual locations in a continuous-space Wright-Malécot model with independent dispersal of all individuals (top) and under our continuous space model incorporating density-dependant fitness (bottom). The clustering seen in the top row is the "Pain in the Torus" described by Felsenstein (1975).

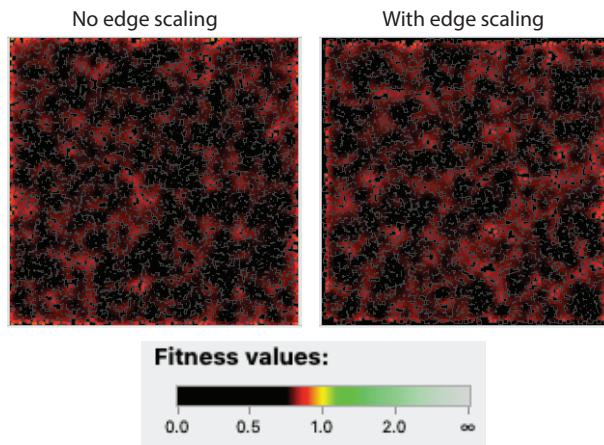


Figure S2 Comparison of individual fitness across the landscape in simulations with (right) and without (left) a decline in fitness approaching range edges. Note the slight excess of high-fitness individuals at edges on the left, which is (partially) counteracted by the scaling procedure.

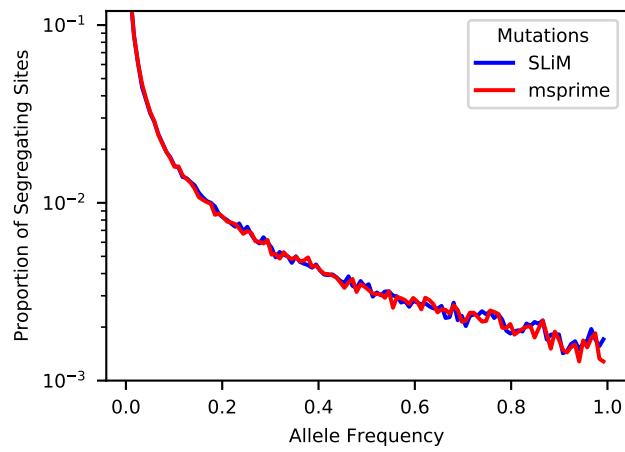


Figure S3 Site frequency spectra from a simulation with neighborhood size = 12.5 when mutations are recorded directly in SLiM (blue line) or applied later in msprime (red line).

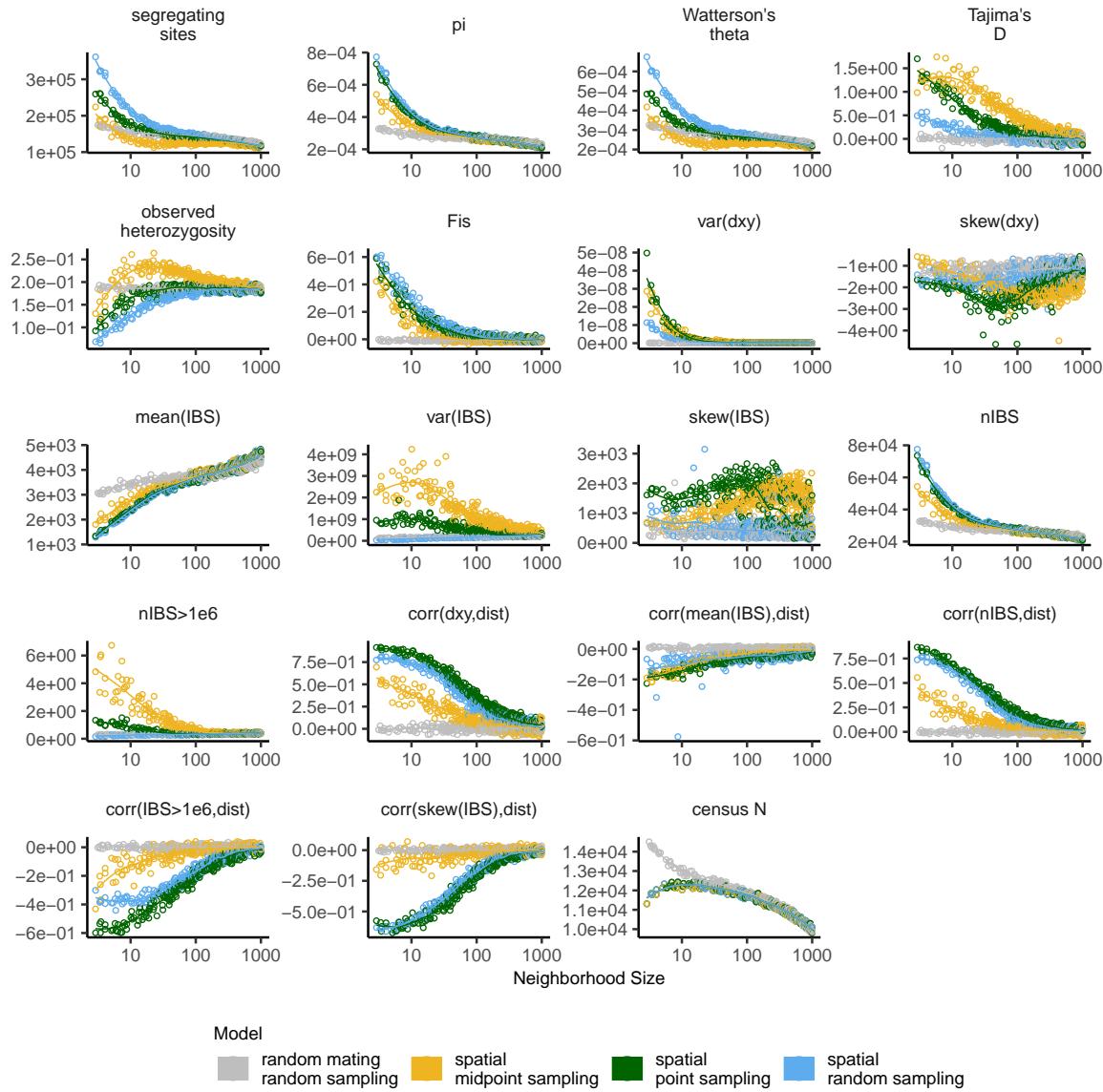


Figure S4 Change in summary statistics by neighborhood size and sampling scheme calculated from simulated sequence data of 60 individuals.

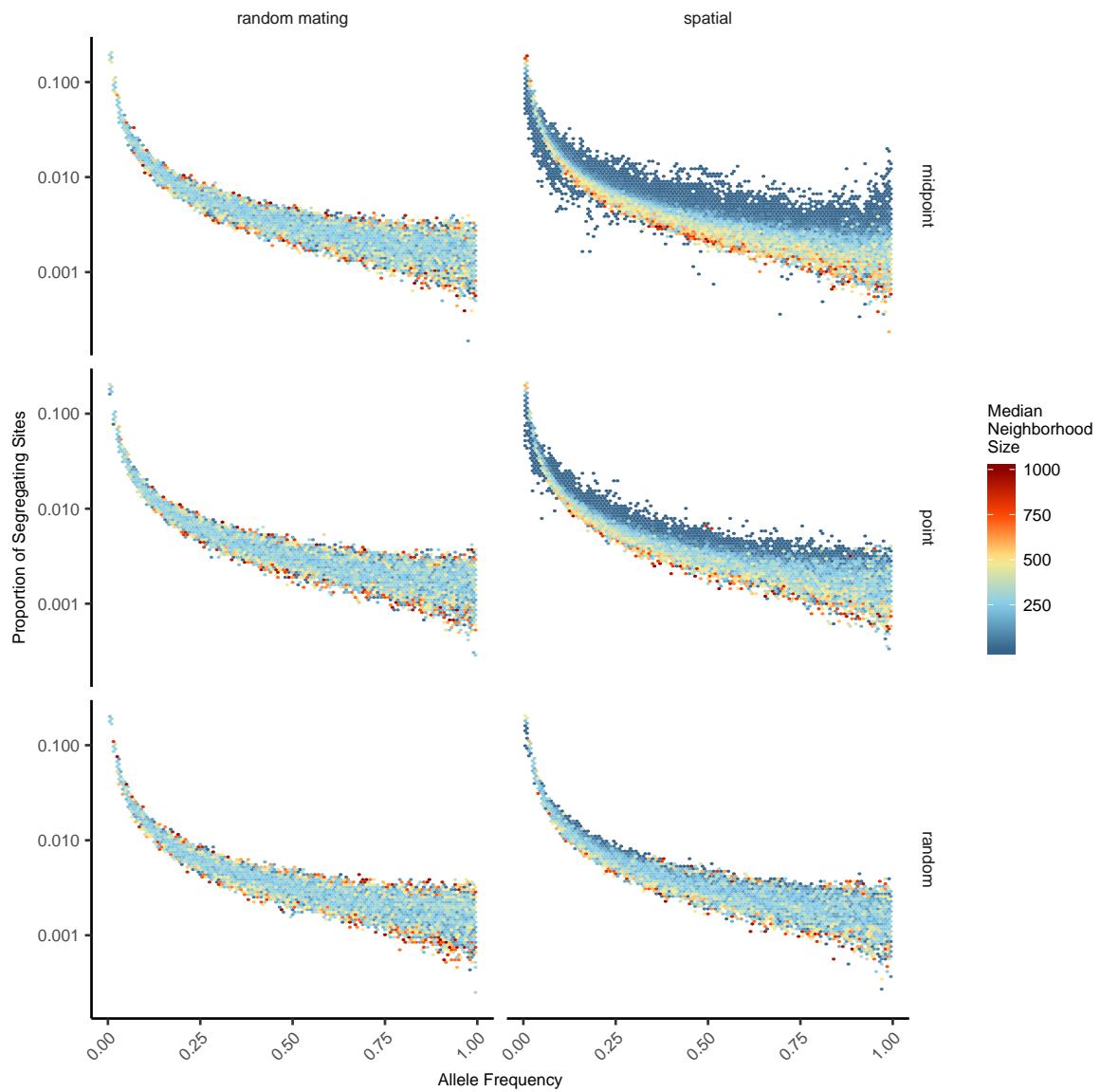


Figure S5 Site frequency spectra for random mating and spatial SLiM models under all sampling schemes.

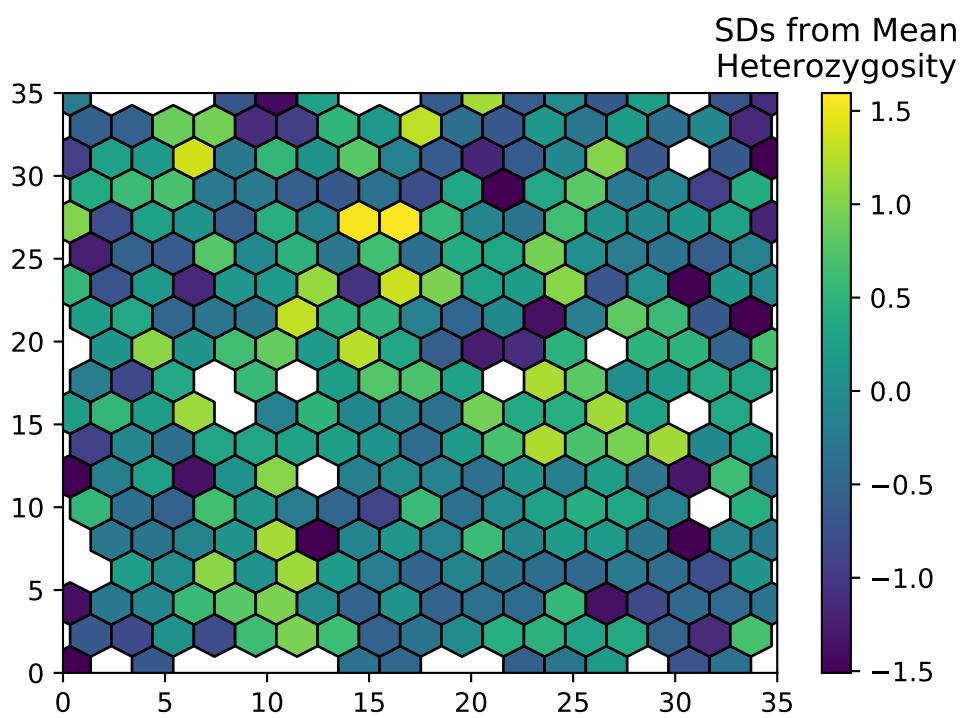


Figure S6 Variation in observed heterozygosity (i.e. proportion of heterozygous individuals) in hexagonal bins across the landscape, estimated from a random sample of 200 individuals from the final generation of a simulation with neighborhood size ≈ 25 . Values were Z-normalized for plotting.

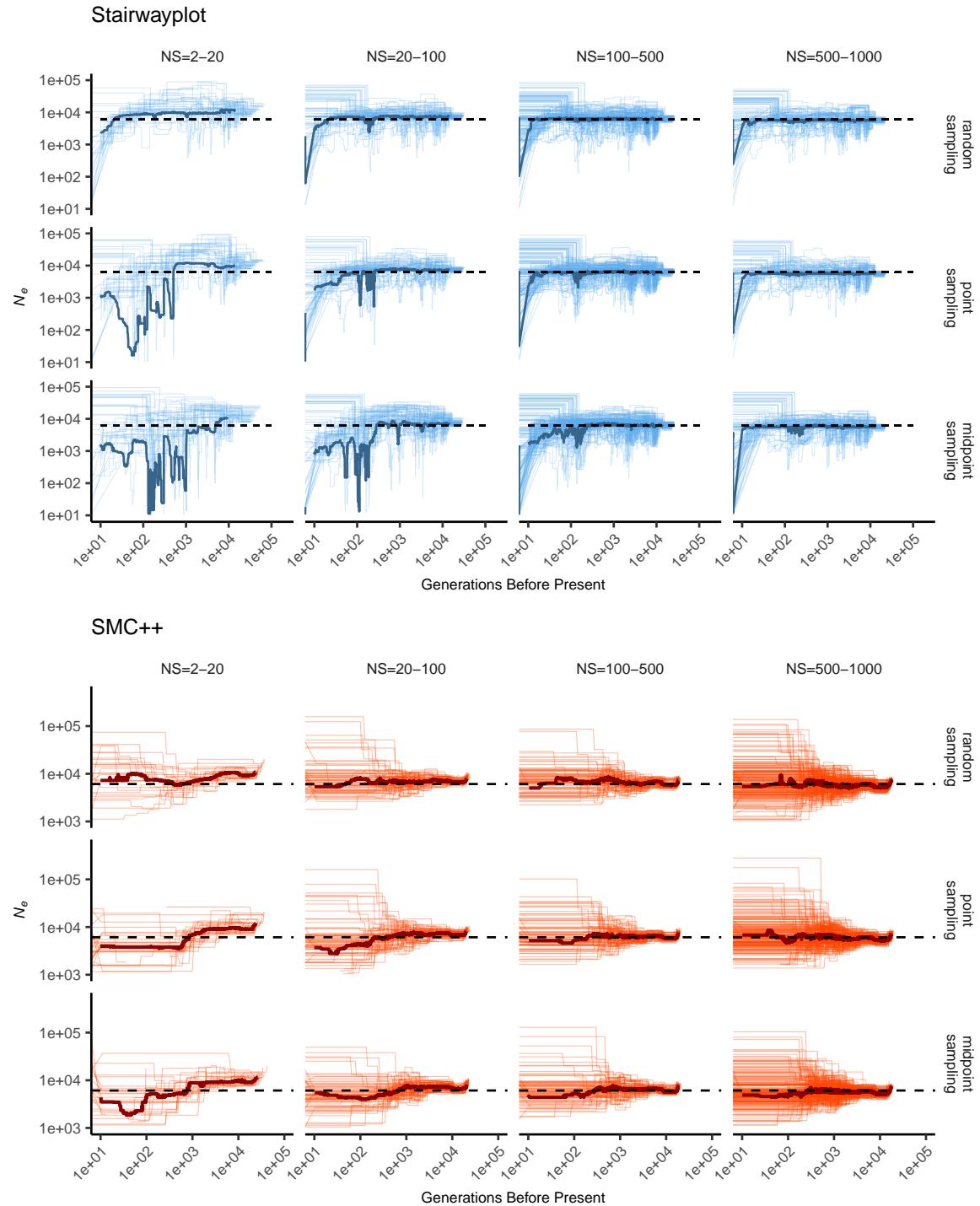


Figure S7 Inferred demographic histories for spatial SLiM simulations, by sampling scheme and neighborhood size (NS) range. Thick lines are rolling medians across all simulations in a bin and thin lines are best fit models for each simulation. Dashed horizontal lines are the average N_e across random-mating SLiM models estimated from θ_π .

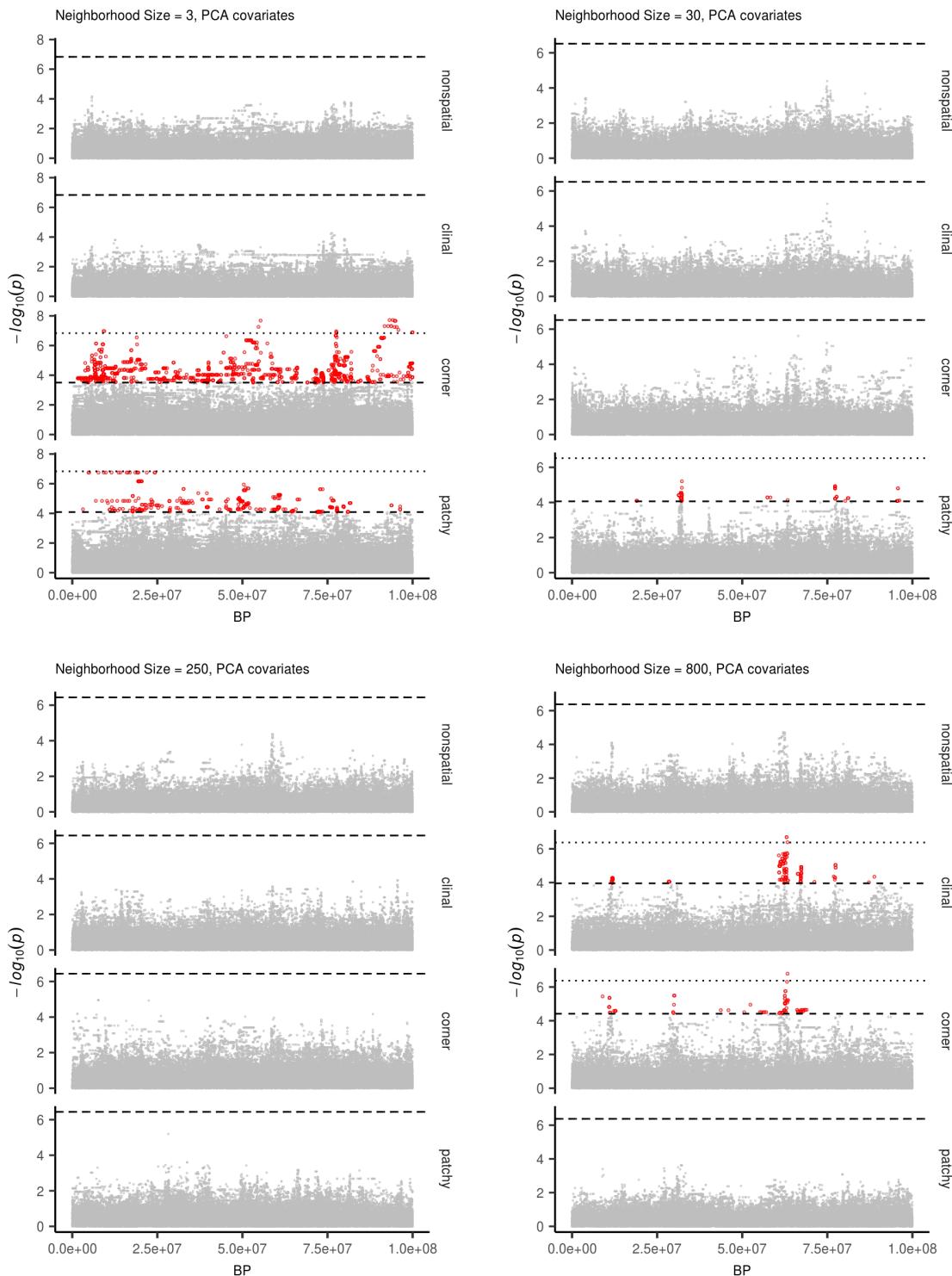


Figure S8 Manhattan plots for a sample of simulations at varying neighborhood sizes. Labels on the right of each plot describes the spatial distribution of environmental factors (described in the methods section of the main text). Points in red are significantly associated with a nongenetic phenotype using a 5% FDR threshold (dashed line). For runs with significant associations the dotted line is a Bonferroni-adjusted cutoff for $p = 0.05$.

Table S1 Summary statistics calculated on simulated genotypes.

Statistic	Description
Θ_{pi}	Mean of the distribution of pairwise genetic differences
Θ_W	Effective population size based on segregating sites
Segregating Sites	Total number of segregating sites in the sample
Tajima's D	Difference in Θ_{pi} and Θ_W over its standard deviation
Observed Heterozygosity	Proportion of heterozygous individuals in the sample
F_{IS}	Wright's inbreeding coefficient $1 - H_e / H_o$
$var(D_{xy})$	Variance in the distribution of pairwise genetic distances
$skew(D_{xy})$	Skew of the distribution of pairwise genetic distances
$mean(IBS)$	Mean of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$var(IBS)$	Variance of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$skew(IBS)$	Skew of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$nIBS$	Mean number of IBS tracts with length > 2bp across all pairs in the sample.
$nIBS > 1e6$	Mean number of IBS tracts over 1×10^6 bp per pair across all pairs in the sample.
$corr(D_{xy}, dist)$	Pearson correlation between genetic distance and $\log_{10}(spatial\ distance)$
$corr(mean(IBS), dist)$	Pearson correlation between the mean of the IBS tract distribution for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(nIBS, dist)$	Pearson correlation between the number of IBS tracts for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(IBS > 1e6, dist)$	Pearson correlation between the number of IBS tracts > 1×10^6 bp for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(skew(IBS), dist)$	Pearson correlation between the skew of the distribution of pairwise haplotype block lengths for each pair of samples and $\log_{10}(spatial\ distance)$

Table S2 Anova and Levene's test p values for differences by sampling strategy. Bolded values are rejected at $\alpha = 0.05$

variable	model	p(equal means)	p(equal variance)
segsites	random mating	0.998190	0.980730
$\Theta\pi$	random mating	0.997750	0.996450
Θ_W	random mating	0.998190	0.980730
Tajima's D	random mating	0.879690	0.188770
observed heterozygosity	random mating	0.531540	0.433230
F_{IS}	random mating	0.474790	0.785730
$mean(D_{xy})$	random mating	0.997770	0.996510
$var(D_{xy})$	random mating	0.283630	0.647240
$skew(D_{xy})$	random mating	0.958320	0.260750
$corr(D_{xy}, dist)$	random mating	0.601980	0.000000
$mean(IBS)$	random mating	0.997960	0.997730
$var(IBS)$	random mating	0.486450	0.399490
$skew(IBS)$	random mating	0.117980	0.069770
$nIBS$	random mating	0.997680	0.996570
$nIBS > 1e6$	random mating	0.834870	0.888730
$corr(mean(IBS), dist)$	random mating	0.073270	0.308420
$corr(IBS > 1e6, dist)$	random mating	0.268440	0.002100
$corr(skew(IBS), dist)$	random mating	0.396920	0.000620
$corr(nIBS, dist)$	random mating	0.581090	0.000000
segsites	spatial	0.000000	0.000000
$\Theta\pi$	spatial	0.026510	0.013440
Θ_W	spatial	0.000000	0.000000
Tajima's D	spatial	0.000000	0.000000
observed heterozygosity	spatial	0.000000	0.000000
F_{IS}	spatial	0.000000	0.000120
$mean(D_{xy})$	spatial	0.025390	0.012910
$var(D_{xy})$	spatial	0.004970	0.006230
$skew(D_{xy})$	spatial	0.000000	0.000000
$corr(D_{xy}, dist)$	spatial	0.000000	0.000000
$mean(IBS)$	spatial	0.272400	0.114250
$var(IBS)$	spatial	0.000000	0.000000
$skew(IBS)$	spatial	0.000000	0.000000
$nIBS$	spatial	0.033920	0.016640
$nIBS > 1e6$	spatial	0.000000	0.000000
$corr(mean(IBS), dist)$	spatial	0.000000	0.590540
$corr(IBS > 1e6, dist)$	spatial	0.000000	0.000000
$corr(skew(IBS), dist)$	spatial	0.000000	0.000000
$corr(nIBS, dist)$	spatial	0.000000	0.000000

Resubmission Cover Letter
Genetics

C. J. Battey,
Peter Ralph,
and Andrew Kern
Monday 18th November, 2019

To the Editor(s) –

We are writing to submit a revised version of our manuscript, “Space is the Place: Effects of Continuous Spatial Structure on Analysis of Population Genetic Data”.

Sincerely,

C. J. Battey, Peter Ralph, and Andrew Kern

Reviewer AE:

The manuscript admirably explores a lot of consequences of isolation-by-distance in the context of a novel model that is easily amenable to forward simulation; however, given that this model may be used in a lot of future studies based on the precedent set here, there is some concern about the model and its support. Reviewers 2 and 3 highlight this in particular (it underlies the main 2 points of reviewer 2's review, and the core of Reviewer 3's comment), and I agree. Whatever can be done to strengthen the standing of this model, and/or connect it to more thoroughly studied models, will be helpful for the manuscript. The concern would be that there are peculiarities of this model that do not generalize well. A new supplemental section or opener to the results section establishing the model more thoroughly would make the strongest response.

*(I would generally cut down the quoted bits like the above to only what's essential, but haven't done that yet.)
(IMPORTANT: don't reorder or delete "points" below - it messes up the automatic numbering!)*

(AE.1) Line 35: Also cite Wilkins and Wakeley, Genetics 2002; Wilkins 2004

Reply: Done. (p. 1, l. 35)

(AE.2) (p. 3, l. 138) “Such models have been used extensively in ecological modeling but rarely in population genetics”: Detailing these previous uses via citations and elaboration may help alleviate the major concern about the provenance of this model and its unique behaviors (see general comments above and R2 and R3 comments).

Reply: Good idea – we have added some historical discussion and a few more citations to this section. (p. 3, l. 138)

(AE.3) (p. 4, l. 185) Please describe computation time needed per replicate

Reply: We have added a figure (Figure 3) and short discussion (p. 9, l. 360) of run times.

(AE.4) (p. 22, l. 776) I read the acknowledgement to the Hearth and Creative Sky Brewing with a sense of familiarity in feeling of gratitude to my own favorite cafes and breweries, but I it's not a great precedent for Acknowledgements to be filled this way. Please cut.

Reply: Good point; we have done this.

(AE.5) Figure 5: Show random-mating expectation

Reply: This is done.

(AE.6) Figure 4A, S5: Perhaps more revealing to show on log-log scale?

Reply: Good suggestion – the SFS in Fig 4A is now on a log-log scale, which shows the slight decrease in low frequency SNPs a little better. We've left Fig S5 as-is.

(AE.7) Figure S3: Caption seems to be missing detail

Reply: Thanks for catching this - we have revised this caption to add details including the simulation parameters. (p. 35, l. 1048)

Reviewer 1:

We thank the reviewer for their very constructive comments. Responses follow below:

This study explores biases arising in population-based inference when 1) real population samples are coming from spatial habitat with various degree of structuring while inference is made assuming random mating population; 2) imperfect sampling in practice that fails to represent full diversity across entire population habitat; 3) phenotypes that vary across geography and create spurious associations with genotypes. While earlier studies explored the effect of strong structure on population genetic inference and GWAS, this work focuses on less extreme scenarios of structuring that arises in populations evolving in continuous habitat. By using non-Wright-Fisher model, authors simulated chromosome-scale samples from populations that evolved in continuous space, and that can model environmental factors to create phenotypes varying over space. As a result, this study identified spatial structuring scenarios (small neighbourhood size 10-100) that coupled with imperfect sampling strategies lead to a biased inference of widely used population genetic statistics (altogether 18 statistics) such as pi (average pairwise sequence differences), heterozygosity (and inbreeding coefficient), and IBS tract sharing. Accordingly, inference of the effective population size history was also strongly affected under these parameter ranges. Finally, the authors use their spatial modelling to demonstrate that typical GWAS with PC-based correction cannot entirely remove spurious signals of genotype-phenotype associations arising from purely environmental factors. Overall, the authors explore an important but often neglected source of bias that can affect inference in many population-based studies (in medical genetics, evolutionary biology and ecology). This study can be of interest to a broader audience of readership, and I have only minor comments to improve clarity and increase accessibility for readers:

(1.1) When neighbourhood size is small (10-100), the mean number of IBS tracts > 2bp (n_{IBS} as in Table S1) is elevated similar to Wright's inbreeding coefficient, but mean of the distribution of pairwise IBS (mean(IBS)) is decreased. What could be the source of this discrepancy? How exactly mean(IBS) was calculated?

Reply: The mean of the IBS tract distribution is calculated by building a (very large) vector of lengths of all IBS tracts between all pairs of individuals and taking the mean over this entire vector. This distribution is quite strongly skewed (i.e. see figure 4 and the skew(IBS) panel of fig S1) such that the highest density is in low values. We explain the depressed mean value by noting that at low neighborhood sizes the coalescence time of distant (in space) individuals is inflated, which allows more time for recombination and mutation to break up IBS tracts. The flip side of this is that there are then more IBS tracts in total across a random sample of individuals, explaining the corresponding inflation in the number of tracts. We also see an increase in the number of very long IBS tracts at low neighborhood sizes, which corresponds to an excess of very recent coalescent events for groups of nearby individuals.

(1.2) The authors use K to denote both carrying capacity (p. 4, l. 163) and population density (p. 4, l. 166). It might be better to use a different notation for these quantities since carrying capacity is fixed while density is an emergent quantity in the non-Wright-Fisher model. Use of K to denote carrying capacity and density is a bit confusing. For example, on (p. 7, l. 319) it is said that 'the "population density" (K) and "mean lifetime" (L) parameters were the same in all simulations'. Here K seems to indicate carrying capacity rather than density? The latter is an emergent quantity and varies across simulation runs?

Reply: We agree that this distinction is worth emphasizing! We've adjusted our language to hopefully remind the reader that K is a parameter that controls population density, rather than being equal to it, at (p. 4, l. 166) and (p. 5, l. 200) and (p. 7, l. 319).

(1.3) Concerning the non-Wright-Fisher model used, it would be helpful to emphasize that some of the parameters are emergent in contrast to Wright-Fisher model. For example, on Page 11, lines 306-308, the author's goal was to look at census size variation and variation in other quantities. This would be better understood if to emphasize that these parameters are emergent properties in the non-Wright-Fisher model used.

Reply: We have added to the text at the beginning of the results to emphasize that this analysis is necessary because these parameters are emergent rather than fixed (p. 7, l. 318).

(1.4) (p. 6, l. 252) Perhaps 'Demographic Inference' might better reflect the content of this section.

Reply: Good suggestion – we have changed the section heading to "demographic inference" (p. 6, l. 252).

(1.5) (p. 7, l. 284) This sentence with 'Gaussian noise with mean zero and standard deviation 10' is confusing since it was mentioned earlier that the modelled phenotype must vary as human height across Europe, and human height varies 2 standard deviations. Only after reading the whole paragraph it becomes clear that 'standard deviation 10' here refers to unit of height. Please consider rephrasing this sentence.

Reply: We have revised this sentence to clarify that we aim to produce a variation in mean phenotype of two standard deviations across the landscape (p. 7, l. 285).

(1.6) (p. 7, l. 306) In the sentence, 'We also examined p values for systemic inflation' I think the authors meant 'systematic inflation'.

Reply: Whoops; thanks. Fixed.

(1.7) Please correct the legend in Figure 2: must be 'spatial model' and 'random mating' model.

Reply: Thanks for catching our confusing legend title placement! We have moved "model" to after "spatial" as suggested.

(1.8) Optional: a dashed line in Figure 2 that shows the total carrying capacity of $50 \times 50 \times 5 = 12500$ would be helpful.

Reply: This is a good suggestion, but we decided to not include this as we don't have straightforward expectations for the other parameters shown.

(1.9) (p. 9, l. 369) The phrase 'affect summaries of variation' is better to replace with 'summaries of genetic variation'.

Reply: Done. (p. 9, l. 369)

(1.10) Please add or correct references to supplementary figures: For example, Figure S2 was probably meant to accompany Figure 3A, while Figure S1 Figure 3B, but references in the text are absent. In fact, the first reference is made to Figure S3 on page 15.

Reply: (check on final pass)

(1.11) There are also several typos and errors in the text. For example, (p. 7, l. 319); (p. 19, l. 654).

Reply: Thank you for noting these – they have been corrected.

Reviewer 2:

Battey et al. use spatially explicit population genetic simulations to analyze the effects of spatial structure on (i) the estimation of key population genetic parameters, in turn used to (ii) make inferences about population history, and on (iii) confounding in genome-wide association studies (GWAS). I liked the paper a lot. It's interesting, well-written and addresses an important question - the effect of spatial population structure on population genetic statistics and inference-and I enjoyed reading it. The most positive aspects were:

1. It nice to actually see spatially explicit simulations and I'm happy that forward simulation is now fast enough that you can do this sort of thing.
2. The paper is very clear and well-written, easy to understand the motivation and most of the details. That's not always the case for this sort of paper.
3. I felt that the section about the effect on GWAS was the most interesting and novel part of the paper and gave me some intuition that I hadn't had before.

I don't have any major criticisms. There were a few aspects that I thought might warrant some additional discussion, and a few specific questions below. The general questions I had after reading it were:

(2.1) To what extent are any of the results dependent on the exact method of simulation. There are a number of choices about the exact details of the simulations (e.g. the way the overlapping generations are handled, the edge effects and, particularly, the form of Equation 1 - see below). It's not so much that these are non-standard (since I don't think there is a standard) and they all sort of make sense heuristically, and I was left wondering whether these sorts of choices actually make a difference. Do the authors have some thoughts/intuition/results about that? Given that the results in Fig. 3 seem quite consistent with expectations, I suspect that on some level it doesn't make much difference but then there are intermediate results like Fig. 2 which seem a bit counter-intuitive and I wonder if those aspects depend on the simulation scheme.

Reply: This is a good point also raised by other reviewers. We have added an appendix section comparing our model to a reverse-time stepping stone simulation on several relevant summary statistics (p. 26, l. 992). We show that for θ_π and several other statistics the stepping stone model approaches the continuous model as the resolution of the landscape increases. We also see some interesting differences that seem to reflect artifacts from discretization. We are also curious how all the other analyses we test would be affected using other simulation schemes, but hope to explore that aspect in future work.

(2.2) Related to the first point, to what extent are the results qualitatively different to those that would be obtained in a stepping-stone model? My interpretation is that they are actually very similar, but I didn't see whether that was explicitly discussed. In some sense, it's still easier to do large simulations in a stepping-stone model so it would be nice to be reassured that that's still ok.

Reply: See above. It is certainly easier and faster, but comes with significant issues either when neighborhood size is lower than the population of a deme, or when demes are small enough that the sample size approaches local N_e , at least in reverse-time simulations.

(2.3) The source of equation (1) is not obvious to me. I sort of see how it makes sense, but a little but more intuition or a brief derivation or an illuminating either in the main text or the supplement, would be helpful.

Reply: Additional discussion of the model is now provided in the appendix (see "Demographic model") (p. 27, l. 1027)

(2.4) The authors use a scaling factor in equation (2) to counteract the increase in fitness of individuals at the edges. Can they provide a figure showing that this is the case. What does "roughly" mean on line 164. Perhaps a heatmap of the fitness of individuals across the grid with and without the scaling factor?

Reply: Good suggestion – we have added a supplemental figure to show the distribution of individual fitness across the landscape with and without our edge-scaling approach. (Note though that we found the effect subtle and hard to visualize well).

(2.5) It would be helpful provide the figure showing that generating mutations during the forward simulations in SLIM is equivalent to applying mutations using msprime on pre-generated trees (line 185)? It sounds like this procedure would underestimate the variance in the number of mutations, since you remove the effect of random generation time. Is this effect small?

Reply: We have added a figure showing sample site-frequency-spectra generated from a subset of simulations run with SLiM mutations, and then using msprime to apply mutations to the same tree sequences with our generation-time scaling approach ???. These approaches yield extremely similar spectra so we believe any error in the approximation is quite small. We also now show variation in simulation times with different mutation and tree sequence recording schemes in Figure 3, and included a brief discussion of this in (p. 9, l. 360)

(2.6) Can the authors provide a bit more intuition behind the patterns of variation seen in generation time, census population size, and variance in the number of offspring with respect to neighborhood size seen in Figure 2? For example, it is not obvious to me why the census population size, for example, should decline systematically with respect to neighborhood size. Presumably this isn't just due to the local demographic stochasticity. Could the authors briefly interpret the observed patterns or cite appropriate literature?

Reply: We have added a brief discussion of the causes of the census size scaling to the relevant results section (p. 8, l. 343).

(2.7) Fig. 7D: I am surprised by the extent to which the observed values of $-\log_{10}(p)$ fall below the $y=x$ line. Particularly in the lower right panel for large neighbourhood sizes. I would expect that to be close to panmictic - why are the P-values underdispersed? That seems like a potential bug, or else something weird is going on.

Reply: We have checked the code to the best of our abilities and did not find a bug causing the underdispersion. It seems to reflect overcorrection in the regression when using PC coordinates as covariates – the PCA is capturing some information about the spatial genetic variation which itself covaries only weakly with the phenotype, and as a result we see anomalously low $-\log_{10}(p)$ when regressing genotype against phenotype.

(2.8) Lines 706-716, It might be worth citing Haworth et al Nature Communications 2019 (<https://doi.org/10.1038/s41467-018-08219-1>) who do the proposed test (GWAS for birth location) in UK Biobank to illustrate the population structure.

Reply: Done - thank you for pointing us to this study. (p. 20, l. 710)

(2.9) The analysis and discussion around the effect of GWAS is focused on PCA correction. Do mixed models help at all?

Reply: We are also very interested to know how mixed models perform here, but think that adding a second GWAS method would make this section too large for the current paper. We have added a note to (p. 19, l. 694) specifically citing mixed models as alternate methods that may perform better.

(2.10) *The github link to the code didn't work for me. I assume it will be made public later, but at this point I can't tell whether the code is available/useable.*

Reply: Apologies, I had forgotten to make the repo public. The link should work for everyone now.

Reviewer 3:

The present study deals with a “hot topic” in spatial population genetics. Most inferential and descriptive methods in statistical spatial population genetic rely on a discrete approximation of space and it is not clear what impact this approximation may have when individuals migrate along a continuum instead. Spatial patterns in sampling is also another major issue which is often simply dismissed, mainly because of the paucity of statistical methods to deal with it. This work touches on these important issues in a timely manner.

Although I was enthusiastic about the topic, I was quite disappointed with the core of the study, i.e., the forward-in-time simulation of populations in continuous space. The field has been struggling with this issue for decades – examples of spectacular failures like the Wright-Malecot model (see Felsenstein’s “pain in the torus” article, 1975) or, more recently, the “mugration” or “discrete trait analysis” model in phylodynamics (see De Maio et al. 2015) have probably mostly harmed our research field – that one cannot make the economy of using a sound probabilistic model for generating geo-referenced genetic data. It does not seem to be the case here unfortunately.

(3.1) *First, the simulation starts with individuals distributed uniformly at random in space. Is there any indication that the three-step algorithm used here maintains this distribution during the course of evolution? If it does not, then is there any stationary regime and how many generations does one need to wait before reaching it? I do appreciate that the competitive interaction term was introduced in order to avoid seeing the “clumping” of individuals that hampers the Wright-Malecot model. Yet, just because there are no such clusters does not mean that the spatial distribution of individuals reaches a stable regime and that the distribution reached, if any, is reasonable from a biological perspective.*

Reply: This is a good point. We evaluated this while developing the simulation using the built-in visualization tools of the SLiM GUI, and have now added a supplementary figure (Figure ??) showing the distribution of individuals in our density-dependant spatial model and a continuous-space Wright-Fisher simulation without density dependance.

Second, the demographic process used here involves birth and death of individuals. Does the population survive asymptotically or, like any birth-death process, eventually dies with probability one? In fact, one needs to know a little about the dynamics of the population size to decide whether the corresponding process is reasonable from a biological standpoint.

(3.2)

Reply: Similarly here we monitored the population size over time while running SLiM. All simulations maintained asymptotic populations, and none of the runs we started ended because populations crashed.

(3.3) *Third, it is not clear what the relationship between the expected lifespan and the probability of survival is. The expected lifespan, L , is first defined as the inverse of the expected number of offspring produced by a parent. The authors also define the probability of survival of a given individual at*

a given point in space, p_i . Hence, the expected lifespan at a point in space (and time) is the mean of a geometric distribution with parameter p_i , i.e., $1/p_i$. Now, it is far from being obvious what the relationship between these two approaches for defining the expected lifespan actually is.

Reply: We have added additional discussion of this point to appendix A (see "Demographic Model") (p. 27, l. 1027)

(3.4) Also, the web page <https://github.com/petrelharp/spaceness> does not seem to exist so that I was not able to experiment with the forward-in-time generator used here unfortunately.

Reply: I apologize – I had forgotten to set the repo to public. It should work now.

(3.5) All in all, more efforts need to be made here in my opinion to show that the forward-in-time simulations generate sensible outcomes. Sensible in terms of the behavior of the population demography at equilibrium (provided such equilibrium indeed exists) along with that of the spatial distribution of individuals. The authors could provide some guarantee of the good behavior of their model as evidenced from simulations using a broad range of parameter values for generating data. Alternatively, they could elect to use the spatial-Lambda-Fleming-Viot model for their simulations, which, in my opinion would seem the most sensible option given that (1) it is possible to run backward-in-time simulations under this model, thereby saving a lot of computation time and (2) it is a well-studied model with good mathematical and biological properties and (3) it is implemented in a publicly available software program (<https://github.com/jeromekelleher/discsim>)

Reply: This is a good point also raised by other reviewers. We have added an Appendix comparing our model to a stepping-stone, rather than the spatial-Lambda-Fleming-Viot, because we think these are the most familiar and widely used class of spatial models. See Appendix 1. We find that many features of our model are well approximated by stepping-stone models, and that for statistics like θ_π the stepping stone model results approach our continuous space model as the number of demes used to describe the landscape increases.

(3.6) Figure 2: I do not understand why the neighborhood size varies to the same extent in the random mating model as it does for the spatial model. For the random mating model, I would have expected the neighborhood size to be equal to the census size since all individuals have the same probability of being a parent of any given offspring. From lines 166->171, it is clear that the spatial model would converge to the random mating model when the mean parent-offspring distance tends to infinity only if we were to ignore the impact of range edges. I am thus wondering whether the variation of neighborhood size one observes in Fig 2 for the random mating model is just a consequence of border effects. If that is the case, then the authors should state it clearly and try to justify it from a biological perspective.

Reply: We have added a brief discussion of the census size scaling to (p. 8, l. 343)

(3.7) Line 729-731: "Many more species occur in a middle range of neighborhood sizes between 100 and 1000 - a range in which spatial processes play a minor role in our analyses [...]" Do the authors think that the spatial processes would still play a minor role when neighborhood sizes exceed 100-1000 if the habitat was larger than that taken in the present simulations? It would also probably be useful to mention that neighborhood sizes given in Table 1 should be compared with extreme caution since the size of the corresponding habitats vary across species. More generally, I suspect that the size of the habitat has a substantial impact on the vast majority of statistics examined in this study. Indeed, the mean parent-offspring distance, which is at the core of the definition of Wright's neighborhood size, is only small or large relative to the size of the habitat.

Reply: This is a good point. Wright's work (Wright 1943) suggests some aspects of genetic variation such as variance in allele frequencies and inbreeding coefficients can be estimated by looking only

at what he would later (Wright 1946) call “neighborhood size”, but certainly other aspects like the number of segregating sites will also depend on total landscape size. We now note on (p. 21, l. 720) that we have evaluated only one landscape size, and have added a sentence to the discussion noting that exploration of these patterns in varying landscape sizes is an important avenue for further research (p. 21, l. 755).

(3.8) Line 753-757: please add a reference to Guindon, Guo and Welch (2016). This study clearly shows that population density and dispersal parameters are identifiable and can indeed be estimated in practice under the spatial Lambda-Fleming-Viot model.

Reply: Done. Thank you for pointing us to this study.

Reviewer 4:

The manuscript by Battey et al explores the consequence of a well-known violation to population genetic models: the fact that populations are spatially structured and mate along a geographical cline, rather than randomly. This topic is important, particularly in light of recent work describing how spatially correlated genetic and environmental impacts can confound some population genetic insights, such as positive selection for height in Europe. The analyses and investigations presented here are thorough and sensible, and my comments are primarily intended to broaden accessibility for this interesting topic.

(4.1) *Introduction. The discussion is very clear, articulating the three primary goals of the project: the impact of failing to model spatial population structure on 1) population genetic summary statistics, 2) inference on demographic history from population genetic data, and 3) impacts on GWAS summary statistics. I found the discussion a bit easier to follow than the introduction and would suggest streamlining and introducing the topic a bit more. Since the paper follows the flow described in the discussion, it might help orient readers by introducing these topics in the same order.*

Reply: Thank you for this suggestion. We have slightly revised the introduction and hope it is now clearer; however since we want to cover a little history and motivation for our continuous model vs stepping-stone approaches in the intro it does have a different flow from the discussion.

(4.2) *I agree that most modern work describes structure as discrete populations connected by migration. However, some methods/studies have explicitly modeled spatial structure, e.g. especially in ecology or using methods like dadi (diffusion approximations). Highlighting some examples of previously identified structure not possible to infer without modeling geography would be helpful to contextualize this work.*

Reply: We have expanded our citations of some of the relevant ecology literature (p. 3, l. 130), which we hope helps to contextualize the study better.

(4.3) *There is some reference to spatial models using grids (e.g. Rousset 1997). Some additional discussion contextualizing more recent methods like EEMS that also construct demes and model migration through divergence between neighboring demes would be helpful and interesting.*

Reply: Good point on EEMS. We have added the most recent EEMS paper to the citations on (p. 1, l. 38), and have added an Appendix giving a more thorough comparison with stepping-stone models.

(4.4) *Demographic modeling. Both approaches tested, stairwayplot and SMC++, are most sensitive to older demographic events, and consequently are very noisy and underestimate effect population sizes, especially in smaller neighborhood sizes. Models that consider haplotype structure are much better suited to this time period. It would be helpful to either 1) discuss the varying time sensitivities*

of different classes of demographic inference methods and how spatial patterns of genetic variation would influence these inferences, or 2) apply a method of this class (many options, e.g. DoRIS, IBDNe, Tracts, Globetrotter, etc) and show how it performs.

Reply: We now discuss haplotype methods in the relevant discussion section (p. 18, l. 627). However though these methods should be more accurate for recent events it is not clear that this will improve performance per se. The dips in recent inferred Ne from stairwayplot are not just prediction noise, but actually reflect an underlying genealogy in which terminal branches are shorter than expected from a constant-size random-mating population (see e.g. figure 4A and 5). The interpretation error is that these short branches are generated by spatial structure rather than changes in population size over time – a point also made in the (Mazet *et al.* 2015) paper we discuss in the introduction and discussion.

(4.5) *GWAS mixed models. To what extent can spatial signals (e.g. corner, patchy) be corrected with mixed models, e.g. with PCs and PC-adjusted GRM as in Conomos *et al.* 2016 using PC-AiR and PC-Relate)? Is patchiness related to dispersal? I'm curious how this relates to the predictive ability of GWAS phenotypes with some spatial association that may or may not be associated with environmental effects.*

Reply: Good question – we are also interested to know how mixed models perform here, but think that to properly test that we would want to change our design to generate phenotypes from simulated genotypes. This would allow us to evaluate false-negatives in addition to false-positives. This is important because, if mixed models do provide stronger control for stratification they are also likely to remove true signals of causal SNPs if those SNPs covary with spatial structure. We now point to these methods explicitly in the discussion (p. 19, l. 694), but think that incorporating that study here would make this paper too long. We also think the PC results are still quite relevant as the method is still seen in many studies.

(4.6) *Code availability. This github link doesn't work, but is important to be able to evaluate for review: <https://github.com/petrelharp/spaceness>*

Reply: Apologies, it was accidentally set to private. The link should work now.

(4.7) *Definitions and interpretations. There are quite a large number of metrics discussed in Figure 3B, and it's a lot to take in. It might be helpful to have a table with a reminder of what the metric is, its interpretation, and how it is computed.*

Reply: We have included a table describing the summary statistics in Figure S1.

(4.8) *Notation: "Offspring disperse a Gaussian-distributed distance away from the parent with mean zero and standard deviation σ in both the x and y coordinates. Each offspring is produced with a mate selected randomly from those within distance 3σ , with probability of choosing a neighbor at distance x proportional to $\exp(-x^2/2\sigma^2)$." I think x may be overloaded here, or I'm confused. Clarify?*

Reply: (Peter and Andy – ideas here?)

(4.9) *When introducing the "spatial model" as opposed to this "random model," the more concrete illustration in Figure 1 hasn't yet been referenced, which makes it harder to follow. It would be helpful to introduce this figure with the model. Additionally, when Figure 1 is introduced, the order is from right to left (random, then point, then midpoint). It would be helpful to rearrange the figure to mirror what's in the text.*

Reply: We have rearranged the figure as suggested.

(4.10) Not sure I follow this example: “Concretely, an individual at position (x, y) in a 50×50 landscape has mean phenotype $100 + 2x/5$.”

Reply:

(4.11) (p. 7, l. 319) Minor typo (through vs though): This occurs because, even through the “population density” (K) and “mean lifetime” (L) parameters...

Reply: Thanks, this sentence has been revised and fixed.

(4.12) Define NS abbreviation in Figure 6.

Reply: Done.