

# Inferring continuous and discrete population genetic structure across space

Gideon S. Bradburd<sup>\*1</sup>, Graham M. Coop<sup>†,2</sup> and Peter L. Ralph<sup>‡,2</sup>

<sup>\*</sup>Ecology, Evolutionary Biology, and Behavior Graduate Group, Department of Integrative Biology, Michigan State University, MI 48824, <sup>†</sup>Center for Population Biology, Department of Evolution and Ecology, University of California, Davis, CA 95616, <sup>‡</sup>Institute of Ecology and Evolution, Departments of Mathematics and Biology, University of Oregon, Eugene, OR 97403

**ABSTRACT** A classic problem in population genetics is the characterization of discrete population structure in the presence of continuous patterns of genetic differentiation. Especially when sampling is discontinuous, the use of clustering or assignment methods may incorrectly ascribe differentiation due to continuous processes (e.g., geographic isolation by distance) to discrete processes, such as geographic, ecological, or reproductive barriers between populations. This reflects a shortcoming of current methods for inferring and visualizing population structure when applied to genetic data deriving from geographically distributed populations. Here, we present a statistical framework for the simultaneous inference of continuous and discrete patterns of population structure. The method estimates ancestry proportions for each sample from a set of two-dimensional population layers, and, within each layer, estimates a rate at which relatedness decays with distance. This thereby explicitly addresses the “clines versus clusters” problem in modeling population genetic variation, and remedies some of the overfitting to which nonspatial models are prone. The method produces useful descriptions of structure in genetic relatedness in situations where separated, geographically distributed populations interact, as after a range expansion or secondary contact. We demonstrate the utility of this approach using simulations and by applying it to empirical datasets of poplars and black bears in North America.

**KEYWORDS** population genetics; isolation by distance; population structure; model-based clustering

## Introduction

A fundamental quandary in the description of biological diversity is the fact that diversity shows both discrete and continuous patterns. For example, reasonable people can disagree about whether two populations are separate species because the process of speciation is usually gradual, and so there is no set point in the continuous divergence of populations when they unambiguously become distinct species. The issue of identifying meaningful biological subunits extends below the species level, as patterns of phenotypic and genetic diversity within and among populations are shaped by continuous migration and drift, as well as by more discrete events, such as rapid expansions, bottlenecks, rare long-distance migration, and separation

by geographic barriers. Both discrete and continuous components are required to accurately describe most species’ patterns of genetic relatedness.

From a practical standpoint, we often need to identify somewhat separable populations from which individuals are sampled (Wright 1949), even while acknowledging continuous processes. Delineating populations is useful for systematics and for informing conservation priorities (Moritz 1994; Waples 1998; Moritz *et al.* 2002). Furthermore, we often need to identify subsets of individuals resulting from reasonably coherent evolutionary histories for downstream analyses to learn about population history and adaptation. Conversely, the substantial information available from continuous, geographic differentiation (e.g., adaptation along a climatic gradient) can be confounded by discrete historical processes (e.g., admixture), requiring methods that can disentangle the two.

There have been many methods proposed to characterize population genetic structure, including generating population phylogenies (Cavalli-Sforza and Piazza 1975; Pickrell and Pritchard 2012), dimensionality-reduction approaches such as principal

doi: 10.1534/genetics.XXX.XXXXXX

Manuscript compiled: Tuesday 7<sup>th</sup> August, 2018

<sup>1</sup>Department of Integrative Biology, Michigan State University, 288 Farm Lane, East Lansing, MI 48824. E-mail: bradburd@msu.edu

<sup>2</sup>These authors contributed equally to this work.

components analysis (Meirmans 2009; Menozzi *et al.* 1978; Novembre and Stephens 2008; Price *et al.* 2006), and model-based clustering approaches (e.g., Pritchard *et al.* 2000; Falush *et al.* 2003; Hubisz *et al.* 2009; Alexander *et al.* 2009; Lawson *et al.* 2012; Raj *et al.* 2014; Huelsenbeck and Andolfatto 2007; Corander *et al.* 2003; Caye *et al.* 2018; Guillot *et al.* 2005). Each of these methods perform best in particular situations, but many can give misleading results when applied to data that show a continuous pattern of differentiation, as that produced by geographic isolation by distance (Wright 1943; Novembre and Stephens 2008; Frantz *et al.* 2009). Here, we will focus on model-based clustering, the most widely used class of approaches for population delineation. (We note that the problem of identifying population clusters is distinct from, though of course related to, the problem of detecting barriers to gene flow between populations, (e.g., Bradburd *et al.* 2013; Petkova *et al.* 2016; Barton 2008; Ringbauer *et al.* 2018).) Existing model-based clustering methods model each individual's genotypes as random draws from a set of underlying, unobserved population clusters, each with a characteristic set of allele frequencies, which are estimated. These underlying frequencies are identical for all individuals assigned to a cluster, regardless of their spatial location. Spatial information has been incorporated into some of these methods, by, for example, placing spatial priors on cluster membership (Guillot *et al.* 2005; Caye *et al.* 2018), but this does not address the underlying issue that these methods assume that allele frequencies are constant in a cluster across the species' range.

*Isolation by distance* refers to a pattern of increasing genetic differentiation with geographic separation, which occurs when geographically restricted dispersal allows genetic drift to build up differentiation between distant locations (Wright 1943). Theoretical work, mostly derived from “stepping-stone” models (Kimura and Weiss 1964; Sawyer 1976; Shiga 1988), gives us some analytical predictions for isolation by distance (Malécot 1969; Slatkin 1985; Epperson 2003), and some theory has been derived for continuous space (Nagylaki 1978; Nagylaki and Barcilon 1988; Barton *et al.* 2002), but substantial work remains to be done (Barton *et al.* 2013). Given the generality of the circumstances that generate a pattern of isolation by distance, it is unsurprising that isolation by distance is very widespread in nature (Meirmans 2012; Sexton *et al.* 2014).

The ubiquity of isolation by distance presents a challenge for models of discrete population structure, as it is frequently difficult to determine whether observed patterns of genetic variation are continuously distributed across a landscape, or instead are partitioned in discrete clusters. This problem can be compounded if sampling is done unevenly or discretely across a population or species' range, and has given rise to a debate in the population genetic literature about how best to describe sets of individuals using continuous clines and discrete clusters (e.g., Serre and Pääbo 2004; Rosenberg *et al.* 2005).

Most existing model-based clustering methods are based on a discrete set of clusters, and so tend to partition continuous variation into spurious clusters with spatially autocorrelated cluster membership (Frantz *et al.* 2009; Meirmans 2012). In analyses of empirical datasets, which often show strong isolation by distance, model-based clustering approaches will therefore tend to overestimate the number of discrete clusters present.

To address this, we set out to develop a model-based clustering method that, when possible, uses isolation by distance to explain observed genetic variation. With an explicit spatial component, discrete population structure need only be invoked

when genetic differentiation in the data deviates significantly from that expected given geographic separation. In this paper, we model genetic variation in genotyped individuals as partitioned within or admixed across a specified number of discrete layers, within each of which relatedness decays as a parametric function of the distance between samples. We also implement a cross-validation approach for comparing and selecting models across different numbers of layers, and we demonstrate the utility of our approach using both simulated and empirical data. The implementation of this method, *conStruct* (for “continuous structure”), is documented and available for general use as an R package at <https://github.com/gbradburd/conStruct>.

## Materials and Methods

**Data** The statistical framework of our approach is conceptually similar to that in Wasser *et al.* (2004), Bradburd *et al.* (2013), and Bradburd *et al.* (2016), although we use a somewhat different summary statistic than in this previous work. The model works with allele frequencies at  $L$  unlinked, bi-allelic single nucleotide polymorphisms (SNPs) genotyped across  $N$  samples. Each “sample” may be a single individual, a collection of individuals from a location, or frequencies estimated from pooled sequencing. From these we compute the *allelic covariance* between samples  $i$  and  $j$ , denoted  $\hat{\Omega}_{i,j}$ , as the expected covariance of distinct individual alleles chosen from each of the two samples at a random locus. More precisely, suppose that we pick a random bi-allelic locus uniformly from the genome, pick a random “reference” allelic state from the two alleles seen at that locus, and, in each sample, draw one random allele, recording  $X_i = 1$  if the allele drawn in sample  $i$  matches the random reference, and  $X_i = 0$  otherwise. Then,

$$\hat{\Omega}_{i,j} = \text{cov}[X_i, X_j]. \quad (1)$$

Because we randomly choose the reference allele, each  $X_i$  behaves marginally as a fair coin — in particular,  $\mathbb{P}\{X_i = 1\} = 1/2$ , so  $\hat{\Omega}_{i,j} = 1/4$  for every  $i$  — all information enters through correlations.

Although we describe this as a covariance between individually drawn alleles,  $\hat{\Omega}_{i,j}$  is in fact also the covariance between the allele frequencies of a randomly chosen allele in samples  $i$  and  $j$ , as long as  $i \neq j$ . The choice of allele does not affect subsequent calculations, and so may be arbitrary, and  $\hat{\Omega}$  can be calculated as (derived in “Allelic covariance and inference”):

$$\hat{\Omega}_{i,j} = \frac{1}{L} \sum_{\ell=1}^L (f_{i,\ell} - 1/2)(f_{j,\ell} - 1/2) \quad \text{for } i \neq j. \quad (2)$$

Here  $f_{i,\ell}$  is the allele frequency in the  $i^{\text{th}}$  sample at locus  $\ell$ . This definition of covariance differs from the usual “genetic covariance” (McVean 2009) in that (a) we do not subtract locus means (to make the statistic insensitive to sample configuration), and (b) we randomly choose a reference allele at each locus (to retain insensitivity to choice of reference allele). As noted in Petkova *et al.* (2016), for  $i \neq j$ , this can also be calculated as  $\Omega_{i,j} = (1 - 2\pi_{i,j})/4$ , where  $\pi_{i,j}$  is the genetic distance calculated from those  $L$  sites, i.e., the proportion of sites at which random samples from  $i$  and  $j$  differ.

**Continuous and discrete differentiation** Clustering approaches to describing genetic variation are useful because population history can often be meaningfully described on a coarse scale by interactions between discrete “populations” whose relationships

are delimited by patterns of glaciation, large-scale migration, mountain ranges, and the like. Here we add a spatial component within each such discrete historical component, which we refer to as a set of “layers” that overlay the modern map. We imagine each layer as a geographically distributed population that extends over the entire sampled range of the populations. As depicted in Figure 1, each sample is composed of a mixture of contributions from each of these layers, with the relative contributions of each layer described by a set of “admixture proportions” (the  $w_i^{(k)}$ ). These layers thus take the place of “clusters” in clustering methods, but we do not adopt this term, as “spatial cluster” suggests a clustering in space, while our layers may contribute to genetic variation across the entire geographic range.

Within each of these layers, allele frequencies have positive covariance at geographically close locations, but this covariance is allowed to decay as geographic distance increases. This pattern of spatial decay reflects how migration between nearby spatial regions homogenizes allele frequency changes that arise locally due to drift, but less effectively homogenizes geographically distant regions, resulting in a continuous pattern of isolation by distance within each layer. There is a fixed amount of covariance between layers, irrespective of spatial location. Within each layer, allele frequencies are expected to change gradually with distance, but observed frequencies can change abruptly at many loci if the proportions of ancestry individuals derive from each layer (the admixture proportions) do so as well.

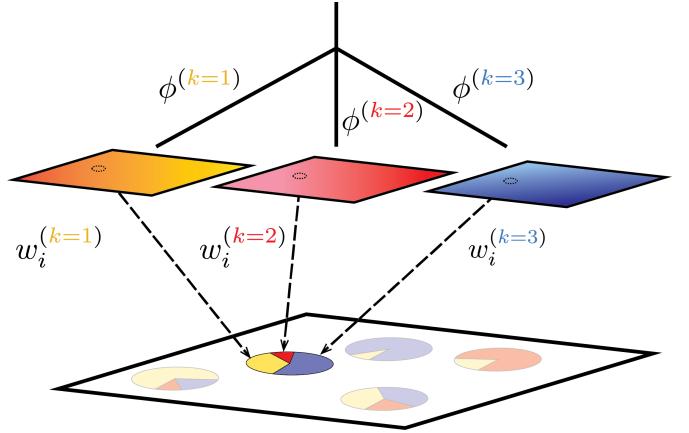
To allow flexibility in the form of the decay of allelic covariance with geographic distance within each layer, we define the covariance within layer  $k$  between samples  $i$  and  $j$  to be:

$$G_{i,j}^{(k)} = \alpha_0^{(k)} \left( \exp \left( -(\alpha_D^{(k)} D_{i,j})^{\alpha_2^{(k)}} \right) \right) + \phi^{(k)} \quad (3)$$

where the superscript  $(k)$  denotes parameters specific to the  $k$ th layer. The quantity  $D_{i,j}$  is the observed geographic distance between samples  $i$  and  $j$ , and the  $\alpha^{(k)}$  parameters control the shape of the decay of covariance with distance in the layer. Our choice of a powered-exponential decay, as parameterized by the  $\alpha$ s, is a flexible and standard choice in spatial statistics (Diggle et al. 1998), and is not chosen to match a particular population genetics model. The  $\phi^{(k)}$  is a parameter that describes the background covariance within the layer. If two samples draw 100% of their ancestry from layer  $k$ , then their covariance under the model is  $G_{i,j}^{(k)}$ ; if they are furthermore geographically very close ( $D_{i,j} = 0$ ) they will have covariance  $\alpha_0^{(k)} + \phi^{(k)}$ . If the geographic distance between them is very large, their covariance will be equal to the background level  $\phi^{(k)}$  within the layer. The “shared drift” parameter  $\phi^{(k)}$  is analogous to the branch length connecting the  $k$ th population to the population ancestral to all modeled layers (see, for example, Patterson et al. 2012; Peter 2016), although they cannot be directly compared because we are modeling the allelic, rather than genetic, covariance. In “Model rationale: drift, admixture, and space” we lay out a simple model of allele frequencies underlying this covariance model.

We then allow samples to draw their ancestry from more than one layer. The admixture proportion of the  $i$ th sample in the  $k$ th layer, denoted  $w_i^{(k)}$ , gives the genome-wide proportion of alleles from sample  $i$  that derive from layer  $k$  (and so  $\sum_{k=1}^K w_i^{(k)} = 1$ ). A visual representation of the method is shown in Fig 1.

We can then describe the covariance between samples  $i$  and  $j$



**Figure 1** Schematic of our method, using  $K = 3$  as an example. Spatial autocorrelation of allele frequencies within each layer is depicted by color gradients, and  $\phi^{(k)}$  denotes the covariance shared by samples with ancestry entirely in the  $k$ th layer. Sampled populations on the landscape are inferred to be admixed between these layers; the  $i$ th sample draws proportion  $w_i^{(k)}$  of its ancestry from layer  $k$ . For convenience, each layer is depicted as a small square, but in fact, each layer exists everywhere in the sampled area, so the small dashed circles on each layer show where the location of the highlighted admixed sample intersects each layer.

across all  $K$  layers,  $\Omega_{i,j}$ , by summing their within-layer spatial covariances ( $G_{i,j}^{(k)}$  in layer  $k$ ), weighted by the relevant admixture proportions.

$$\Omega_{i,j} = \gamma + \sum_{k=1}^K w_i^{(k)} w_j^{(k)} G_{i,j}^{(k)} + \delta_{i,j} \eta_i. \quad (4)$$

In this equation,  $w_i^{(k)} w_j^{(k)}$  is the proportion of alleles that *both* sample  $i$  and sample  $j$  have inherited from layer  $k$ .

In addition to the admixture-weighted sum of the within-layer spatial covariances, this function contains two terms,  $\gamma$  and  $\delta_{i,j} \eta_i$ . The first,  $\gamma$ , describes the global allelic covariance between all samples, and arises because all samples share an ancestral mean allele frequency at each locus, which generates a base-line covariance. In the final term,  $\delta_{i,j}$  is an indicator variable that takes a value of 1 when  $i$  equals  $j$  and 0 otherwise, and  $\eta_i$  adds variance specific to sample  $i$ . This term on the diagonal of the parametric covariance matrix captures processes shaping variance within the sampled deme, such as inbreeding and the sampling process.

**Likelihood and inference** If the allele frequency deviations at each locus were independent between loci and multivariate normally distributed across populations, their allelic covariance  $\widehat{\Omega}$  would be Wishart distributed with degrees of freedom equal to  $L$ , the number of loci genotyped. We use this as a convenient approximation to the true distribution described above, and so define the likelihood of the allelic covariance to be

$$P(\widehat{\Omega} | \Omega) = \mathcal{W}(L\widehat{\Omega} | \Omega, L), \quad (5)$$

where  $\mathcal{W}$  is the Wishart likelihood function. Statistical nonindependence between loci (linkage disequilibrium) will decrease

the effective number of degrees of freedom. One possible solution, which we have not yet found necessary to implement, would be to estimate an *effective* number of loci by introducing a parameter to modify the given degrees of freedom and thereby informally model linkage between loci (e.g., Petkova *et al.* 2016).

We estimate the values of the parameters of the model using a Bayesian approach. Acknowledging the dependence of the parametric covariance matrix  $\Omega$  on its constituent parameters  $w, \alpha, \phi, \eta, \gamma$  and on the (observed) geographic distances  $D$  with the notation  $\Omega(w, \alpha, \phi, \eta, \gamma, D)$ , we denote the posterior probability density of the parameters as:

$$P(w, \alpha, \phi, \eta, \gamma | \hat{\Omega}) \propto P(\hat{\Omega} | \Omega(w, \alpha, \phi, \eta, \gamma, D)) \times P(w)P(\alpha)P(\phi)P(\eta)P(\gamma), \quad (6)$$

where  $P(w)$ ,  $P(\alpha)$ ,  $P(\phi)$ ,  $P(\eta)$ , and  $P(\gamma)$ , are prior distributions. All parameters are given (half-)Gaussian priors except for  $\alpha_2$ , which is uniform on  $(0, 2)$ , and  $w$ , for which we use an independent Dirichlet of dimension  $K$  for each sample (see Table 1 for specifics). Parameters are independent between layers. We use Hamiltonian Monte Carlo as implemented in STAN (Carpenter 2015; Hoffman and Gelman 2014; Stan Development Team 2015, 2016) to estimate the posterior distribution on the parameters. Our R package, `conStruct` (for “continuous structure”), functions as a wrapper around this inference machinery.

**Relationship of this model to nonspatial structure models** A nice feature of our approach is that the model described in Eq. (4) contains a nonspatial assignment model as a special case (see “Models, parameters, and priors” for a more in-depth discussion). By setting  $\alpha_0^{(k)}$  to zero for all  $k$ , we obtain a nonspatial model in which each cluster has its own allele frequency at each SNP, and individuals draw a proportion of their ancestry from each cluster. This model is very similar to that of STRUCTURE (Pritchard *et al.* 2000) and related models (e.g., Alexander *et al.* 2009); the main difference is that our likelihood assumes that allele frequencies are normally distributed around their expectations, while the standard assignment methods assume that the error is binomially distributed (Engelhardt and Stephens 2010). (We make this approximation for the substantial advantages in computational speed.) The second difference is that, in the original STRUCTURE model, allele frequencies at each locus are independently drawn for each cluster (Pritchard *et al.* 2000), while in `conStruct`’s nonspatial model, it is more natural to envision each cluster’s allele frequency as being drifted away from a single, global allele frequency. This makes our model more closely related to the “F-model” prior for allele frequencies of (Falush *et al.* 2003). These differences in the underlying model could in principle result in different behavior, but below we show that the nonspatial model indeed produces similar results to ADMIXTURE, and use this fact to compare the fit of the different models — spatial vs. nonspatial, across different values of  $K$  — by comparing their performance in a common framework.

**Choice of layer number and cross-validation** There are a number of reasons why there is no true (or right) number of layers for real datasets, discussed further in the Discussion. However, it is still important to assess whether additional layers (larger  $K$ ) meaningfully model patterns in the data or merely explain spurious variation introduced by noise — in other words, whether additional model complexity provides significant explanatory power. Toward that end, we have implemented a method for

statistically comparing `conStruct` results across different values of  $K$  and between the spatial and nonspatial models.

Several approaches have been used as model choice criteria for the number of discrete clusters in population genetic data, including: comparisons of the likelihood of the data across different values of  $K$ , with various criteria on how to choose a single value (e.g., Evanno *et al.* 2005), or with information theoretic penalizations such as AIC or BIC (e.g., Alexander *et al.* 2009); comparisons of the marginal likelihood, generated either via various approximations (e.g., Pritchard *et al.* 2000) or via thermodynamic integration (Verity and Nichols 2016); or inference using a Dirichlet process prior (Huelsenbeck and Andolfatto 2007). See Verity and Nichols (2016) for a discussion of these approaches and comparison between several methods.

We use cross-validation (similar in spirit to Alexander and Lange 2011) to attack this problem. To do this, we use a “training” partition of the data (in practice, a random 90% subset of the loci) to estimate the posterior distribution of the parameters, and then calculate the log-likelihood of the remaining “testing” loci, averaged over the posterior. Prediction accuracy of a particular value of  $K$  is then measured using this log-likelihood, averaged over a number of independent data partitions. The best model is judged to be the simplest one with significantly better predictive accuracy than others (see “Cross validation procedure” for more on our cross-validation procedure). In general, larger values of  $K$  allow the model more flexibility, and thus increases the likelihood of the training partition, but this improvement in the likelihood will plateau (or even peak), as above a certain  $K$  the model only fits noise specific to the training data rather than generalizable patterns. At any value of  $K$ , support for the spatial model over the nonspatial model means that isolation by distance is likely a feature of the data.

Cross-validation provides a valuable summary of how much explanatory power is added by spatial structure within each layer, and each additional layer. However, we remind users that “statistical significance does not imply real-world significance”, and so small but statistically significant differences between models should not be relied on too strongly.

Another way to describe the practical significance of additional layers is to calculate each layer’s relative contribution to total covariance, and to choose a value of  $K$  where all layers have a contribution above some cutoff (e.g. 0.1%). The Dirichlet prior on admixture proportions is quite harsh against intermediate admixture values (see Table 1), encouraging the model to “not use” unnecessary layers if they are present in the model, so that they will have a low contribution to overall covariance.

To calculate layer contributions, we use the following alternative description of our covariance model: the genomes of any pair of individuals *agree* with some background probability at a locus, but this probability of agreement is increased on any segment of genome that both have inherited from the same layer (the amount it increases depends on how far apart they are geographically and on the decay of isolation by distance). We use this characterization to quantify the relative contributions of each layer, by computing the average contribution to increased probability of agreement as described in “Calculating layer contributions”. This layer contribution is similar to the “ancestry contribution” proposed by Raj *et al.* (2014). However, each of our layers can induce a different amount of covariance between samples embedded in them, so we take that into account when calculating each layer’s contribution to the whole.

## Data Availability

The method `conStruct` is implemented as an R package, and is available for installation at <https://github.com/gbradbard/conStruct>. Scripts for generating and analyzing all simulated and empirical datasets, as well as the datasets themselves, are also available at the same site, and additionally have been archived at Data Dryad (doi:10.5061/dryad.5qj7h09).

## Results and Discussion

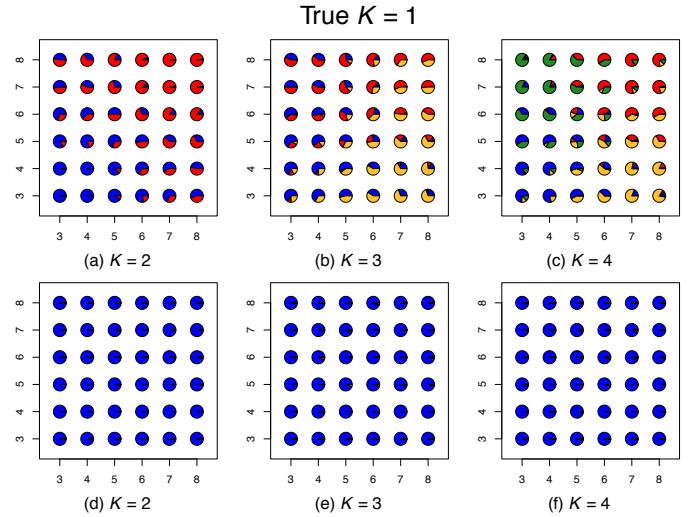
### Simulations

To test the method, we first generated data using the coalescent simulator `ms` (Hudson 2002). In each simulation, we split a single ancestral population into  $K$  subpopulations  $\tau_s$  units of coalescent time in the past, and at time  $\tau_e$  in the past, each of these discrete populations instantaneously colonized a separate  $6 \times 6$  square lattice of demes. Migration on each lattice was to nearest neighbors (eight neighbors, including diagonals). Finally, at time  $\tau_a$  in the past, we collapsed those  $K$  discrete layers into a single grid of demes, choosing various amounts of admixture from these different layers (see Fig 10), with randomly distributed but spatially autocorrelated admixture proportions. See “Simulation details” for more details, including parameter values used. We simulated datasets using  $K = 1, 2$ , and  $3$  layers; in each simulation we sampled 10,000 unlinked loci from each of 20 haploid individuals from every deme. We then ran both spatial and nonspatial `conStruct` analyses on each simulated dataset with  $K$  between 1 and 7, and compared predictive performance of the models using cross-validation with 10 replicates. For comparison, we also analyzed each simulated dataset using ADMIXTURE (Alexander *et al.* 2009) with  $K$  between 1 and 7, and compared models using ADMIXTURE’s cross-validation procedure with 50 folds.

With these simulations, spatial `conStruct` does not create spurious discrete groupings when there are none: Figures 2 and S1-S3 show that subsequent layers beyond the number used for simulation are unused. When data simulated with  $K = 1$  are analyzed with  $K > 1$ , the additional layers contribute very little to any population. Even when the spatial model is run with  $K = 7$ , the inferred admixture proportions are nearly identical to those estimated under the true value of  $K$  for each simulation. Moreover, the method infers the true admixture proportions with high accuracy, tight precision, and good coverage (Figs S4-S5).

In contrast, the nonspatial model describes geographic variation using gradients of admixture between increasingly many discrete clusters to better approximate the continuous, spatial patterns of relatedness (Figs 2 and S6-S8). The ADMIXTURE results are qualitatively similar, as shown in Figs S9-S11. Each nonspatial cluster is genetically more similar within itself than it is to other clusters, but we know that these boundaries are arbitrary, because the data were simulated without them.

The spatial model’s better fit is reflected by increased predictive accuracy: as shown in Fig 3, across all models and choices of  $K$ , the spatial model is correctly preferred over the nonspatial model. As desired, predictive accuracy of the spatial model increases until the true value of  $K$ , and then plateaus or declines (Figs 3 and S12-S14). Predictive accuracy of the nonspatial model increases as subsequent clusters are added up to  $K = 7$  (the largest number tested), although gains are greatest as layers below the true number are added. The same holds true for the ADMIXTURE cross-validation results, in which models that



**Figure 2** Results for data simulated using  $K = 1$ , showing maps of admixture proportions estimated using the nonspatial `conStruct` model for  $K = 2$  through 4 (top row) and the spatial `conStruct` model for  $K = 2$  through 4 (bottom row). As there is only a single layer in the simulation, no populations should be admixed, which is accurately depicted by the spatial model (second row), while the nonspatial model creates spurious clusters (first row).

have the largest number of clusters are preferred over all other models, as shown in Fig 3 (vermillion diamonds), and, in more detail, in Fig S15.

The unimportance of spurious layers can be seen in plots of layer contributions (Figs 4, S16-S17). In the spatial analyses, once we pass the true  $K$ , subsequent layers add little in terms of (co)variance explained; in contrast, additional clusters in the nonspatial analyses continue to contribute substantially

### Empirical Applications

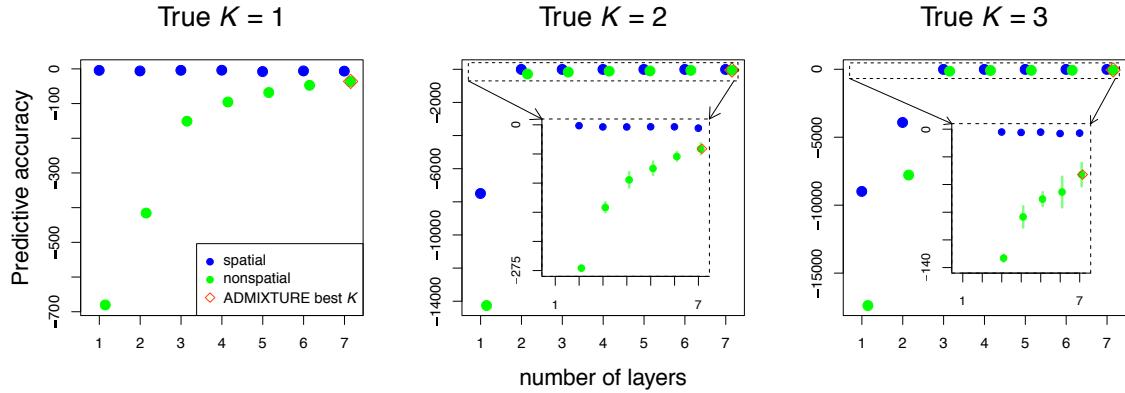
To further demonstrate the utility of this method, we also applied `conStruct` to empirical population genomic data from two systems: a contact zone between two poplar species in northwestern North America, and a large North American sample of black bears.

### Poplars

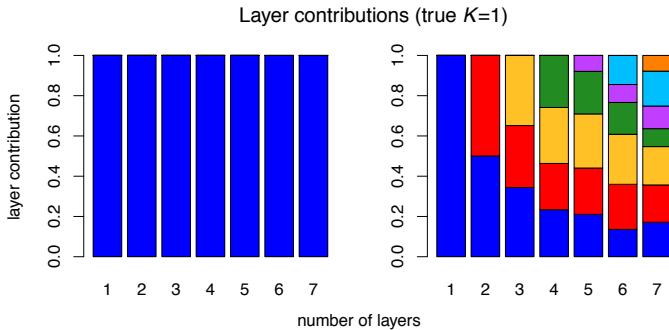
**Study system and questions** Trees in the genus *Populus* (poplars, aspens, and cottonwoods) are distributed throughout the Northern Hemisphere; species in the genus regularly co-occur and, where they do, they frequently hybridize (Eckenwalder 1984; Cronk 2005).

*Populus trichocarpa*, the black cottonwood, and *Populus balsamifera*, the balsam poplar, have a broad zone of overlap in the Pacific Northwest, where they are hypothesized to hybridize (Geraldes *et al.* 2014; Suarez-Gonzalez *et al.* 2016). Both species are sampled over a large geographic region, and show spatial patterns of genetic and phenotypic variation (Slavov *et al.* 2012; McKown *et al.* 2014), making the system well-suited for application of our method. We organize the results of our analyses around the following questions:

1. To what degree has hybridization blurred the boundaries between *trichocarpa* and *balsamifera*? (As an extreme case, does



**Figure 3** Cross-validation results for data simulated under  $K = 1$ ,  $K = 2$ , and  $K = 3$ , comparing the spatial and nonspatial conSTRUCT models (in blue and green, respectively) run with  $K = 1$  through 7, with 10 cross-validation replicates. The inset plots zoom in on cross-validation results outlined in the dotted boxes. The spatial model shows better model fit at every value of  $K$ . The vermilion diamond indicates the value of  $K$  selected on the basis of lowest cross-validation error among ADMIXTURE models. In all simulations, the preferred ADMIXTURE model was that with the largest number of clusters.



**Figure 4** Results for data simulated using  $K = 1$ , showing layer/cluster contributions (i.e., how much each layer/cluster contributes to total covariance), from conSTRUCT runs using  $K = 1$  through 7 for the spatial model (left), and the nonspatial model (right). In each run of the spatial model, a single layer explained nearly all the covariance (additional bars are present but not visible).

genetic differentiation support these as separate species, as opposed to a single cline of ancestry?)

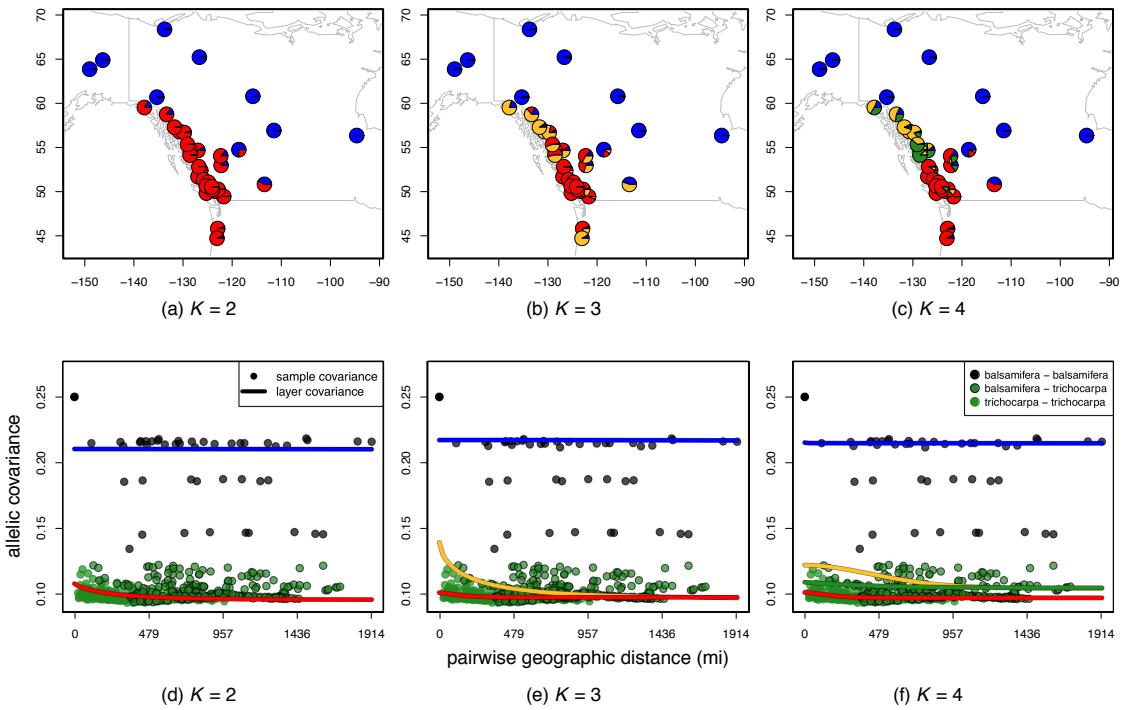
2. Does the only significant boundary of population structure fall along the species boundary (if any), or is there substructuring within species?
3. Does the strength of isolation by distance differ between inferred layers? This may indicate, e.g., different speeds of postglacial expansion or primary modes of dispersal.

**Data and analyses** We use data from Geraldes *et al.* (2014), consisting of 434 individuals sampled from 35 drainages genotyped at just over 33,000 loci (map of the sampling shown in Fig S18). The number of individuals per drainage ranged between 1 and 50, with most sampling concentrated on *trichocarpa* drainages. The data were generated using an Infinium 34K array designed for *trichocarpa* (Geraldes *et al.* 2013), and showed a strong pattern of bias in allelic dropout (the majority of missing data were from drainages with only *Populus balsamifera* individuals). To ameliorate some of the problems that arise when there is a strong bias in which data are missing, we dropped loci for which any data

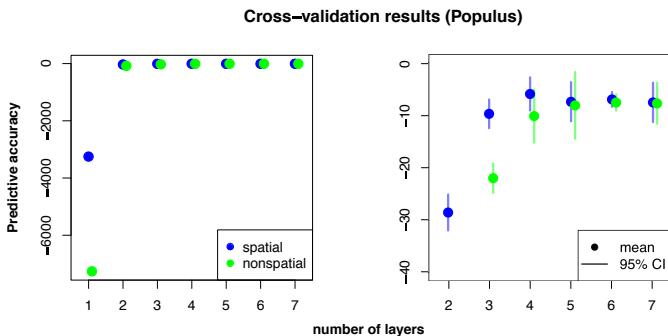
were missing, resulting in just over 20,200 loci retained for analysis. We then analyzed these data, grouped by drainage, using both the spatial and nonspatial conSTRUCT models with  $K = 1$  through 7, and compared these models using cross-validation with 10 replicates. The results of all these analyses are shown in Figs 5 and 6, as well as Figs S19-S23 in the “Supplementary Materials”. For comparison, we also ran ADMIXTURE (Alexander *et al.* 2009) with  $K = 1$  through 7, using 50-fold cross-validation to compare model performance (Figs S24-S25).

**Results from conSTRUCT** All models with  $K > 1$  assigned the majority of each of the two species to distinct layers, with some populations drawing ancestry from multiple layers. Based on cross-validation results, we view the  $K = 3$  spatial model as a sufficient description of the data, with additional structure of uncertain significance. This provides strong support for discrete population structure between the two species, with some admixture, rather than a single, continuous cline of ancestry. At all values of  $K > 1$ , discrete population structure was mostly partitioned along species lines; at values of  $K$  above 2, further discrete substructure was inferred within the *P. trichocarpa* samples, with no substructure within *balsamifera*. There was also strong support for isolation by distance in the dataset, but most of this signal seems to derive from the *P. trichocarpa* samples: as seen in Figs 5d-f and S21, there is almost no isolation by distance within the *balsamifera* layer ( $\alpha_D \approx 0$ ). Both points are in agreement with previous work (Keller *et al.* 2010), which found low diversity within the region’s *balsamifera*, probably as the result of a recent postglacial expansion.

A consistent split between layers within *trichocarpa* fell along the “no-cottonwood belt,” a region along the central coast of British Columbia in which black cottonwood is absent (the break between yellow and red, for  $K \geq 3$ ). The no-cottonwood belt is hypothesized to divide the species’ distribution into northern and southern groups, which, in a provenance test, were experimentally shown to display differences in ecologically relevant phenotypes (e.g., pathogen resistance, Xie *et al.* 2009; Chang-Yi *et al.* 2012). At higher values of  $K$ , drainages at the southern tip of *trichocarpa* sampling begin to split out into their own layers, perhaps due to introgression from the southern neighbors *P. angustifolia* or *fremontii* (Zhou and Holliday 2012; Geraldes *et al.* 2014).



**Figure 5** Maps of admixture proportions estimated for the *Populus* dataset using the spatial `conStruct` model for  $K = 2$  through 4 (a-c), as well as the corresponding layer-specific covariance curves estimated under each model (d-f).



**Figure 6** Cross-validation results for *Populus* dataset comparing the spatial and nonspatial `conStruct` models run with  $K = 1$  through 7 with 10 cross-validation replicates. The first panel in each row shows all results; the second panel zooms in on the results from analyses run with  $K = 2$  through 7.

**Comparison to ADMIXTURE** Both nonspatial `conStruct` and ADMIXTURE displayed the successive partitioning of space and the clines of admixture seen in the simulation results. The details of each were somewhat different (Fig S20 vs. Fig S24), and also differed across the replicate analyses. These differences between runs and methods may be due to noise in the different inference algorithms employed, multi-modality in the likelihood surfaces, or to model details (e.g., the priors used in nonspatial `conStruct`, or the fact that ADMIXTURE is modeling each allele's frequency in each cluster, rather than the covariance across all alleles). However, overall, the behavior of both methods was quite similar: each recovered the *trichocarpa/balsamifera* split with the first two clusters modeled, then, with higher values of  $K$ , used subsequent clusters to subdivide the *trichocarpa* samples into geographically restricted foci of cluster membership. Both

nonspatial `conStruct` and ADMIXTURE strongly favored the most cluster-rich model (Figs 6 and S25). In contrast, the spatial `conStruct` model clearly did not favor the model with the highest value of  $K$ , and appears to describe patterns of isolation by distance across the *trichocarpa* range quite well.

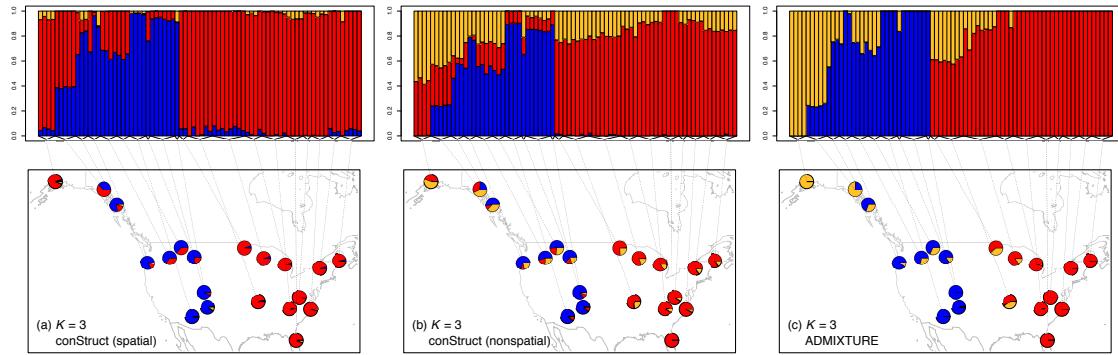
### Black bears

**Study system and questions** The American black bear, *Ursus americanus*, is endemic to North America and has a broad distribution across the continent. During the last glacial maximum, black bears were confined to isolated glacial refugia, from which they subsequently expanded to occupy their current range (Wooding and Ward 1997; Byun *et al.* 1997; Stone and Cook 2000; Puckett *et al.* 2015), likely leading to both continuous and discrete patterns of genetic structure. We organize our results around the following questions:

1. How many distinct populations are reflected in modern patterns of genetic variation?
2. How strong is isolation by distance within each inferred group?

Distinct populations likely represent different glacial refugia, and differing strengths of isolation by distance might indicate different levels of habitat connectivity, dispersal behavior, or different postglacial histories.

**Data and analyses** We use data from Puckett *et al.* (2015), consisting of 95 individuals sampled across the United States and on the West coast of Canada, genotyped at just under 22,000 bi-allelic loci. The distribution of missing data across these individuals was uneven, with a few individuals representing most of the missing data, so we removed individuals with greater than 4% missing data, resulting in a final dataset of 78 individuals.



**Figure 7** Maps of admixture proportions estimated for the black bear dataset using the spatial `conStruct` model (left), the nonspatial `conStruct` model (center), and `ADMIXTURE` (right) for  $K = 3$ . Pies show mean admixture results across individuals within their diameter, and the admixture results for all individuals included within each group are shown in the plot above.

We then analyzed these data, treating individuals as the unit of analysis, using both the spatial and nonspatial `conStruct` models with a  $K$  of between 1 and 7, and compared these models using cross-validation with 10 replicates. We also ran `ADMIXTURE` (Alexander *et al.* 2009) on the same dataset, using  $K = 1$  through 7, and comparing models using `ADMIXTURE`'s cross-validation procedure with 50 data fold subsets. The results of these analyses are shown in Figs 7–9, as well as in Figs S26–S31 in the “Supplementary Materials”.

**Results from `conStruct`** The results partition the sampled bears into two main groups (shown in Fig 7a): one (red) to the east of the Rocky Mountains, which also occurs in Alaska, the other primarily west of the Rockies (blue). The disjointed range of the red layer likely reflects the fact that Canada was not sampled, and so the red layer may extend through the intervening (unsampled) northern Great Plains and Canadian Shield, with the blue layer presumably then stretching up into British Columbia.

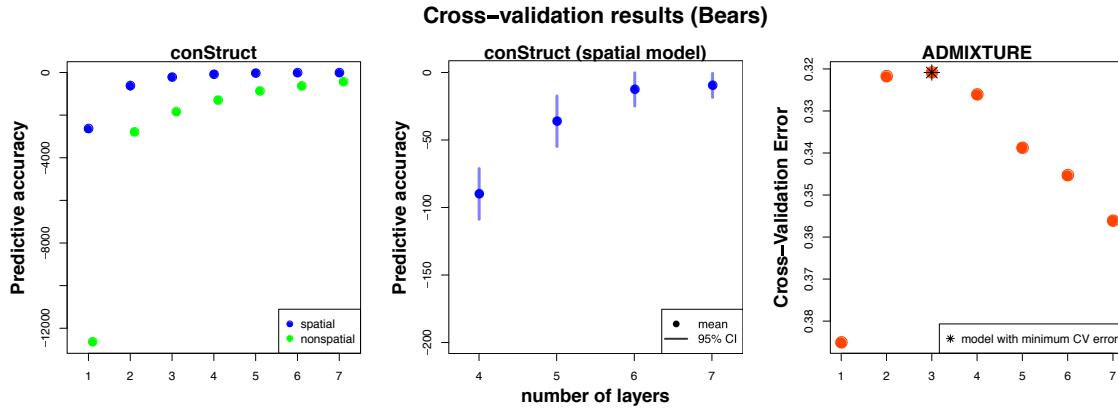
The spatial models have strong statistical support up until around  $K = 5$  or 6 (Fig 8), but additional spatial layers beyond  $K = 2$  contribute little to total covariance (Fig 9). The locations of admixed individuals are consistent with a scenario of postglacial expansion from two refugia, one in the American Southwest and one in the American Southeast, meeting near the Northwest coast of North America and the Cascade Range. However, lack of any samples from Canada and Mexico, and lack of denser sampling across northern North America, make more detailed interpretations untrustworthy. The spatial covariance functions estimated in layers beyond the first two take very large values over small spatial lags, but decay sharply after that. This feature, combined with the overall amounts and spatial patterns of ancestry in those layers, suggests that these layers are describing processes that shape genetic variation at local scales, such as inbreeding, which affects covariance between individuals within each location, but has limited impact on covariance between locations.

**Comparison to `ADMIXTURE`** Results from the nonspatial model and from the `ADMIXTURE` analyses clearly exhibit the tendency of nonspatial clustering algorithms to describe continuous spatial patterns of divergence using gradients of admixture between clusters. For example, in Fig 7b, the third cluster (in gold) exhibits a clear East-West gradient that overlays the discrete structure between the Southwest cluster and the Southeast. The results from `ADMIXTURE` are not identical to those obtained using the nonspatial `conStruct` model, but they do show the

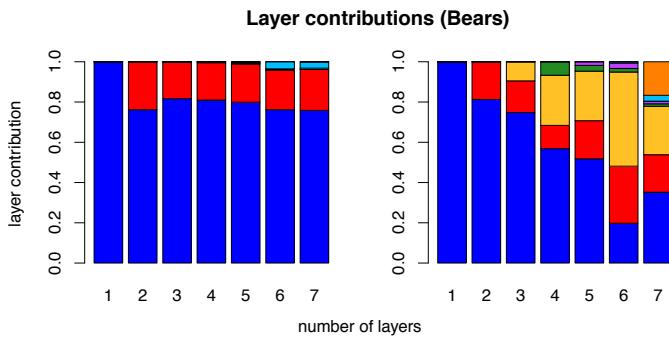
same tendency: e.g., at  $K = 3$  — the preferred model from the cross-validation analysis shown in Fig S31 — `ADMIXTURE` splits the westernmost Alaskan samples out of the cluster with the eastern samples, and at  $K = 4$ , it subdivides the eastern cluster into two geographically partitioned groups (Fig S30). Interestingly, for the nonspatial model implemented in `ADMIXTURE`, the preferred model has a smaller  $K$  ( $K = 3$ ) than that of the spatial models with best cross-validation performance in `conStruct` ( $K = 5$  or 6). This discrepancy likely stems from the different features introduced in layers beyond  $K = 2$  in the two models: `conStruct` uses small contributions of new layers to model very local drift, while `ADMIXTURE` moves to geographically finer subdivisions.

Even at  $K = 3$ , `ADMIXTURE` invokes clusters to describe what seems to be a continuous spatial pattern of genetic variation, which `conStruct` describes using only two spatial layers. The third cluster in the `ADMIXTURE` analysis at  $K = 3$  (shown in gold in Fig S30b), shows strong spatial autocorrelation in admixture proportions, as would be expected if it is describing continuous spatial differentiation. The allelic covariances plotted against distance (see Fig S32) provide more information on `ADMIXTURE`'s lack of fit: covariance between Eastern bears falls off gradually rather than abruptly with distance, indicating a residual pattern best explained by isolation by distance within layers. In addition, the covariance between bears assigned to `ADMIXTURE`'s gold and red layers (the furthest Northwestern and Eastern bears, respectively) appears to be a natural extension of the decay of covariance with distance, falling to only slightly lower values than covariances between other widely separated pairs of Eastern sampling locations.

Across all values of  $K$  for which we ran `conStruct`, we see strong support for the spatial model over the nonspatial model (Fig 8). This pattern may resolve a discrepancy between our results and previous analyses that split Alaskan and British Columbian bears out into their own cluster with an inferred Beringian glacial refugium (Byun *et al.* 1997; Stone and Cook 2000; Puckett *et al.* 2015). Our model, which explicitly incorporates a spatial decay of relatedness, allows somewhat genetically differentiated individuals that are sampled far from one another to belong to the same layer, instead of splitting these individuals out into successive clusters (e.g., Fig S26d vs S27d).



**Figure 8** Cross-validation results for the black bear dataset, comparing spatial and nonspatial conSTRUCT models, as well as ADMIXTURE, all run with  $K = 1$  through 7, with 10 cross-validation replicates for the conSTRUCT analyses and 50 data-fold subsets for the ADMIXTURE analyses. The first panel in each row shows results from spatial and nonspatial conSTRUCT models; the second panel zooms in on the results from the spatial analyses run with  $K = 4$  through 7, and the third panel shows the results for ADMIXTURE. Note that the admixture plot shows cross-validation error (rather than predictive accuracy), and that the y-axis has therefore been flipped for ease of comparison to the conSTRUCT results.



**Figure 9** Layer/cluster contributions (i.e., how much total covariance is contributed by each layer/cluster), for all layers estimated in runs using  $K = 1$  through 7 for the spatial model (left), and for all clusters using the nonspatial conSTRUCT model (right). For each value of  $K$  along the x-axis, there are an equal number of contributions plotted. Colors are consistent with Fig 7.

## Discussion

In this paper, we have presented a statistical framework, conSTRUCT, for simultaneously modeling continuous and discrete patterns of population structure. By employing the sensible default assumption that relatedness ought to decay with geographic distance, even within a population, we avoid erroneously ascribing population differentiation to discrete population clusters. To aid comparison between models, we present a cross-validation approach as well as a way to describe the contribution of each spatial layer to the model (but caution against overly strict interpretation of either).

The method performs well on simulated data: we accurately infer the admixture proportions used to simulate the data and accurately pick the simulating model as the best model using our cross-validation procedure. Two empirical applications of conSTRUCT to samples of North American poplars and black bears yield reasonable results, and demonstrate that, by acknowledging isolation by distance, real datasets can be better described using fewer layers.

The proposed method combines the utility of model-based clustering algorithms with a model of isolation by distance. We anticipate that conSTRUCT will be useful for identifying populations and determining ancestral origins of sampled individuals, especially when the populations exhibit geographic patterns of relatedness.

**Comparison to nonspatial model-based clustering** Above, we showed that (a) the nonspatial conSTRUCT model recapitulates results of other, commonly-used nonspatial clustering methods, and (b) conSTRUCT can concisely capture spatial structure, which is common within populations. Given this, when should methods without spatial capability be used? One advantage these have over conSTRUCT is speed when the number of samples is large. Although conSTRUCT's computation time is independent of the number of loci included in the dataset (after the initial calculation of the allelic covariance), it currently scales poorly with number of samples. The computationally limiting step is the inversion of the covariance matrix, which scales more than quadratically with the number of samples, whereas computation time for, e.g., STRUCTURE, scales linearly with number of samples.

For a relatively small number of samples, conSTRUCT can be much faster than existing nonspatial Bayesian clustering methods. On a desktop machine, using a single 4.2 GHz Intel Core i7 processor, an analysis of the black bear dataset (78 samples, 21,000 loci) running conSTRUCT's spatial model with 4 layers for 5,000 MCMC iterations (which was more than sufficient for convergence) took 2.8 hours. For almost any size dataset, the maximum likelihood algorithm implemented in ADMIXTURE is quite a bit faster than conSTRUCT: running ADMIXTURE on the bear dataset over all values of  $K$  from 1 to 7, including 50-fold cross-validation for each value of  $K$ , took only 6.6 minutes. It should also be noted that there may be situations when the binomial-based model underlying ADMIXTURE performs better than our Gaussian-based model, e.g., when clusters differ at only a few strongly differentiated loci, although we have not investigated this possibility.

**Choosing the “best” number of layers** Although we recognize the utility of choosing a single, “best” value of  $K$ , and using

only that analysis to communicate results, we emphasize that the choice of best  $K$  is always relative to the data in hand and the questions to be answered. From a statistical perspective, unless the data were generated under the model itself, the support for larger values of  $K$  is likely to increase with increasing amounts of data. In the limit of infinite data, the best value of  $K$  may be the number of samples included in the dataset (Patterson *et al.* 2006).

From a biological perspective, it is important to stress that patterns of relatedness between individuals and populations are shaped by complex spatial and hierarchical processes. All individuals within a species are related to one another in some way, and summarizing those patterns of relatedness with a single value of  $K$  may be reductive or misleading. We therefore encourage users to perform analyses across different values of  $K$  and observe which layers split out at what levels (this is conceptually similar to taking successively shallower cross-sections of the population phylogeny), and also to take the results of the proposed cross-validation procedure with a large grain of salt. Calculating layer contributions may also be a useful heuristic, as it can reveal layers with statistical support but small biological import.

Although we believe our model adds spatial realism to the groups used by clustering methods, it is important to note that the layers detected by our method do not necessarily correspond to distinct, ancestral populations; nor does a non-zero admixture proportion indicate that admixture (i.e., gene flow) must have occurred. Both groupings and admixture proportions should be viewed as hypotheses that should be subject to further testing (for an in-depth discussion of these points, see Falush *et al.* 2016).

**Implications for management and conservation** Because isolation by distance is common, a likely result of applying conStruct to existing data is that populations previously identified as distinct using nonspatial clustering methods may be grouped into the same layer. This “lumping” might better reflect the demographic history of these populations, but may not contradict the genetic distinctness implied by the nonspatial clustering. This genetic distinctness – rather than shared history – may be more relevant for management decisions and conservation policy, which are often predicated on the identification of discrete “management units” identified using genetic data (Moritz 1994; Waples 1998; Moritz *et al.* 2002).

It is therefore important to stress that individuals sampled from the same conStruct layer may be quite genetically diverged from one another, perhaps especially at loci underlying adaptive traits, and that a conStruct layer may still contain multiple distinct management units worthy of independent protections. For instance, although both the Alaskan and Eastern Black Bears draw most of their ancestry from the same conStruct layer, they are separated by a great distance, and may therefore differ substantially from each other (although less than from the Western bears, as measured by average covariance). Alternatively, the inclusion of multiple management units into a single conStruct layer may occur if these populations are currently (or were recently) exchanging migrants, and thus might emphasize the importance of maintaining habitat corridors, or of implementing an integrated conservation plan across a geographic region.

**Allelic or genetic covariance?** The choice of allelic covariance, rather than genetic covariance, was motivated by the fact that it is less affected by sample configuration – the genetic covari-

ance is calculated after subtracting the mean from the entire sample, which is more strongly affected by densely sampled locations. Genetic covariance is also often computed after first dividing each frequency by  $\sqrt{p(1-p)}$ , where  $p$  is the global allele frequency, with the aim of equalizing variances across loci. Our definition does not do this, and so is less affected by low-frequency alleles. Both of these changes led to better performance on test data. However, note that allelic covariance is more affected by singleton sites than the standard genetic covariance, so it may be advisable to filter these prior to analysis if they are likely to contain a large percentage of errors (Linck and Battey 2017).

**Caveats and considerations** There are a few important caveats to consider in the interpretation of conStruct results. First, we have modeled allelic covariance within a layer as a spatial process. Although there is flexibility built into the model about the shape of that covariance, inference may be misleading if the sampling geography departs radically from the way the sampled organisms disperse (or have dispersed) on their landscape. For example, if we were to run a conStruct analysis using geographic distances between sampled individuals of greenish warblers (Irwin *et al.* 2001) or *Ensatina* salamanders (Wake and Schneider 1998) — two canonical examples of ring species — we might get misleading results. This is because distance between locations on either side of the species’ distributions (across the Tibetan plateau and the Central Valley, respectively) is not representative of the path traversed in the coalescent of a pair of alleles sampled at those locations.

A second caveat is that, in some instances, membership in the same layer may not mean that samples are particularly related. If covariance within a layer decays sharply with distance, and the layer-specific relatedness parameter  $\phi^{(k)}$  is low, individuals separated by a large spatial distance may be in the same layer but have very low pairwise relatedness. It is possible that this is happening in Fig S19. At  $K = 3$ , the southernmost populations of *P. trichocarpa* are in the gold layer, whose other neighbors are to the north, with an intervening group of populations in the red layer, and at  $K = 5$ , those southernmost samples split out and become their own layer. Furthermore, note that in this case  $\alpha_0^{(k)}$  and  $\phi^{(k)}$  are confounded, so differences in  $\phi$  between layers should not be overinterpreted. Again, we encourage users to run analyses across multiple values of  $K$  and to examine the spatial covariance functions within layers when interpreting results.

**Extensions and future directions** There are several ways in which the model described in this paper might be extended or improved. For example, we currently assume that all layers within a model are equally unrelated (a star population phylogeny, although the branches can have different lengths thanks to the  $\phi^{(k)}$  parameter), similar to the F-model of (Falush *et al.* 2003). However, we could extend the existing model by implementing a relatedness structure between the layers by, for example, estimating a population phylogeny between them (e.g., Pickrell and Pritchard 2012).

In addition, here we have assumed that samples have known geographic coordinates, and that they draw ancestry from layers only at those sampled locations. A natural extension would be to attempt to “geo-locate” the ancestry of samples without geographic coordinates (Wasser *et al.* 2004). We could also imagine letting samples draw ancestry from other geographic coordinates, as we have done in a previous approach (Bradburd *et al.*

2016) to model long distance dispersal. We could even allow entire layers to bud off of a particular location on another layer. This would enable more explicit modeling of range expansion or domestication, in which a set of individuals are thought to have ancestry that originated from a particular geographic location embedded in a larger pattern of isolation by distance.

A final direction would be to model relatedness within a layer as a spatiotemporal process, in which covariance decays both with distance in space and in time. As the number of genotyped historical or ancient samples increases, it is becoming possible to ask whether there is genetic continuity at a point in space across time, or whether populations are being replaced (Lazaridis *et al.* 2014; Haak *et al.* 2015; Slatkin and Racimo 2016; Nielsen *et al.* 2017; Schraiber 2017; Joseph and Pe'er 2018). However, we expect allele frequencies to change through time in a population, even without replacement, simply due to drift. Therefore, a natural way to test for population replacement is to estimate the rates at which relatedness within a layer decays with time in the same way we do in the current model with space, in which case a change in discrete population structure across space is comparable to population replacement across time.

## Appendices

### Model rationale: drift, admixture, and space

Here we sketch a simple model of allele frequencies and their covariances, to justify the form given in the main text.

**Drift** We first provide a simple model of allele frequencies within a layer. Imagine a sample  $i$  that draws all of its ancestry from layer  $k$ . The allele frequency in sample  $i$  at locus  $\ell$ , denoted  $F_{i,\ell}$ , can be written as the sum

$$F_{i,\ell} = \epsilon_\ell + \Delta_\ell^{(k)} + \Delta_\ell^{(k,i)} + \Delta_\ell^{(i)}. \quad (7)$$

The first term is the ancestral allele frequency  $\epsilon_\ell$  shared by all samples; the second is the deviation from that ancestral frequency due to drift in the ancestral population of the  $k$ th layer, which is shared by all samples within the layer. The third term is the deviation of the  $i$ th sample away from the  $k$ th layer mean due to the spatial process of drift and migration within the layer. The final term is the deviation specific to the  $i$ th sample, which captures drift not shared by all samples at the population level (i.e., subpopulation-specific drift due to, e.g., inbreeding). We will assume that these four deviations are all uncorrelated with each other.

If we have two samples  $i$  and  $j$  drawn from layer  $k$ , their covariance across loci will be

$$\text{Var}(\epsilon) + \text{Var}(\Delta^{(k)}) + \text{Cov}(\Delta^{(k,i)}, \Delta^{(k,j)}) + \delta_{i=j} \text{Var}(\Delta^{(i)}), \quad (8)$$

where the quantity  $\delta_{i=j}$  is an indicator variable that equals 1 when  $i$  is equal to  $j$  and 0 otherwise, as in Eq. (4).

**Admixture** The model above describes the simple case in which samples draw 100% of their ancestry from only a single layer each. To accommodate admixture between layers, we model sampled genomes as drawn from allele frequencies that are weighted averages of the local frequencies in each layer from which they draw ancestry. The weights,  $w_i^{(k)}$ , describe the “admixture proportion” of sample  $i$  in layer  $k$ . These can be interpreted as the proportion of the genome in the  $i$ th sample that came from the  $k$ th layer (or the probability that an allele at a locus is drawn from layer  $k$ ), so that  $\sum_{k=1}^K w_i^{(k)} = 1$  for each  $i$ . The

allele frequency in the  $i$ th sample at the  $\ell$ th locus can therefore be written as:

$$F_{i,\ell} = \epsilon_\ell + \sum_K w_i^{(k)} (\Delta_\ell^{(k)} + \Delta_\ell^{(k,i)}) + \Delta_\ell^{(i)}, \quad (9)$$

and so the covariance between  $i$  and  $j$  across loci is

$$\begin{aligned} \Omega_{i,j} = & \text{Var}(\epsilon) + \sum_{k=1}^K w_i^{(k)} w_j^{(k)} (\text{Var}(\Delta^{(k)}) + \text{Cov}(\Delta^{(k,i)}, \Delta^{(k,j)})) \\ & + \delta_{i=j} \text{Var}(\Delta_i). \end{aligned} \quad (10)$$

**Space** Under our nonspatial model, we assume that  $\text{Cov}(\Delta_\ell^{(k,i)}, \Delta_\ell^{(k,j)}) = 0$ , so that the only additional covariance between  $i$  and  $j$  (above that induced by a shared ancestral frequency at each locus) is due to the drift in the ancestral population of their layer (the variance of which is  $\phi^{(k)}$ ). Under our spatial model we assume that some of the covariance in allele frequencies between  $i$  and  $j$  decays as a function of the geographic distance between the pair,  $D_{i,j}$ , so that

$$\text{Cov}(\Delta_\ell^{(k,i)}, \Delta_\ell^{(k,j)}) = \alpha_0^{(k)} \times \left( \exp\left(-(\alpha_D^{(k)} D_{i,j})^{\alpha_2^{(k)}}\right) \right). \quad (11)$$

We note that this form is chosen for flexibility and convenience, and not because it matches any explicit population genetic model of isolation by distance.

### Allelic covariance and inference

Here we go into further detail about both the allelic covariance we model and the modeling framework we use.

**Allelic covariance** To see why equations (1) and (2) for the allelic covariance are equivalent, pick a random locus and let  $A$  and  $B$  be randomly drawn alleles at that locus from populations  $i$  and  $j$  respectively. Suppose these are each coded as ‘0’ or ‘1’ (where ‘0’ denotes a reference allele), but we randomly “flip” this coding, so that we let  $X = A$  and  $Y = B$  with probability 1/2, but otherwise we let  $X = 1 - A$  and  $Y = 1 - B$ . These are  $X_i$  and  $X_j$  in equation (1), so that  $\hat{\Omega}_{i,j} = \text{cov}[X, Y]$ . The random allele flipping makes the value of  $\hat{\Omega}$  independent of the choice of reference allele. By conditioning on the flip, and using the fact that  $\mathbb{E}[X] = \mathbb{E}[Y] = 1/2$ , Eq. (2) comes from the observation that

$$\text{cov}[X, Y] = \mathbb{E}[(A - 1/2)(B - 1/2)]. \quad (12)$$

Thanks to averaging over choice of alleles, the within-population allelic variance in sample  $i$ ,  $\hat{\Omega}_{i,i}$ , is the variance of a series of Bernoulli(1/2) draws across loci, and therefore  $\hat{\Omega}_{i,i} = 1/4$  for every sample  $i$ . Averaging over choice of reference allele therefore removes some information about factors acting within populations that might otherwise leave signatures in the genetic covariance, such as population size, extent of inbreeding, and history of bottlenecks. However, as our model is focused on modeling covariances between samples as the outcome of some spatial process, we count this a minor loss.

**Likelihood** If allele frequency deviations are well approximated by a Gaussian, their sample allelic covariance is a sufficient statistic, so that calculating the likelihood of their sample allelic covariance is the same as calculating the probability of the frequency data up to a constant. We can therefore model the covariance

of the sample allele frequencies,  $\hat{\Omega}$ , as a draw from a Wishart distribution with degrees of freedom equal to the number of loci  $L$  across which the sample allelic covariance is calculated:

$$\hat{\Omega} \sim \mathcal{W}(L\Omega, L) \quad (13)$$

where  $\mathcal{W}$  is the Wishart likelihood function.

A benefit of directly modeling the sample allelic covariance is that, after the initial calculation of the sample covariance matrix, the computation time of the likelihood is not a function of the number of loci, so inference can be done using whole genome data.

### Models, parameters, and priors

Here we discuss the different models implemented in this paper and give the priors we place on model parameters.

**Spatial versus nonspatial** In this paper, we discuss two types of models, spatial and nonspatial, each of which can be implemented with different numbers of layers/clusters. The spatial model is parameterized as in Eq. (10), and the nonspatial model is a special case of the spatial model with all  $\alpha$  parameters set to 0. The nonspatial model therefore has  $3K$  fewer parameters than the spatial model, because there are three  $\alpha$  parameters that describe the continuous differentiation effect of distance in each layer.

**Single layer** Each of these models can be run with a single layer ( $K = 1$ ), in which case the layer-specific covariance parameter  $\phi^{(k)}$  and the global covariance parameter  $\gamma$  become redundant. The single-layer model is therefore a special case of the multi-layer model, in which we set  $\phi$  to zero. For the spatial model, the single-layer parametric covariance is:

$$\Omega_{i,j} = \text{Var}(\epsilon) + \alpha_0^{(k)} \times \left( \exp \left( -(\alpha_D^{(k)} D_{i,j}) \alpha_2^{(k)} \right) \right) + \delta_{i=j} \text{Var}(\Delta_i), \quad (14)$$

and for the nonspatial model, it is:

$$\Omega_{i,j} = \text{Var}(\epsilon) + \delta_{i=j} \text{Var}(\Delta^{(i)}). \quad (15)$$

**Priors** We use a Bayesian approach to parameter inference. A table of all parameters, their descriptions, and their priors is given in Table 1.

### Cross validation procedure

We employ a Monte Carlo cross-validation approach for model comparison (Picard and Cook 1984). This procedure generates a mean predictive accuracy for each model and each value of  $K$ , as well as a confidence interval around that mean, which can then be used for model comparison or selection. Briefly, we follow the following procedure:

1. For each of  $X$  replicates:
  - (a) partition the allele frequency data into a 90% “training” partition ( $F_1^X$ ) and a 10% “testing” partition ( $F_2^X$ )
  - (b) run our inference procedure using the training partition to estimate model parameters  $\theta_{mk}$  for  $2K$  models:
    - i.  $m$ : the spatial and the nonspatial model
    - ii.  $k$ : the number of layers/clusters 1 through  $K$

Parameter	Description	Prior
$\gamma$	global covariance due to shared ancestral frequency	$\gamma \sim \mathcal{N}(\mu = \text{Var}(\vec{f}), \sigma = 0.5)$
$\alpha_0^{(k)}$	controls the sill of the covariance matrix in layer $k$	$\alpha_0^{(k)} \sim \mathcal{N}(\mu = 0, \sigma = 1)$
$\alpha_D^{(k)}$	controls the rate of the decay of covariance with distance in layer $k$	$\alpha_D^{(k)} \sim \mathcal{N}(\mu = 0, \sigma = 1)$
$\alpha_2^{(k)}$	controls the shape of the decay of covariance with distance in layer $k$	$\alpha_2^{(k)} \sim U(0, 2)$
$\eta_i$	the nugget in population $i$ (population specific drift parameter)	$\eta_i \sim \mathcal{N}(\mu = 0, \sigma = 1)$
$\phi^{(k)}$	layer-specific shared drift in layer $k$	$\phi^{(k)} \sim \mathcal{N}(\mu = 0, \sigma = 1)$
$w_i$	admixture proportions sample $i$ draws across $K$ layers	$w_i \sim \text{Dir}(\alpha_1 \dots \alpha_K = 0.1)$

**Table 1** List of parameters used in the conStruct model, along with their descriptions and priors. The mean of the Normal prior on  $\gamma$ ,  $\text{Var}(\vec{f})$ , is the variance of the sample mean allele frequencies across loci.

- (c) calculate the mean log likelihood of the testing data partition over the posterior distribution of training-estimated parameters for each model ( $\bar{\mathcal{L}}(F_2^X \mid \theta_{mk})$ , henceforth  $\bar{\mathcal{L}}_{xmk}$ )
- (d) generate standardized mean log likelihoods,  $\mathcal{Z}_{xmk}$ , across all models run on this data partition:
  - i. identify the highest mean log likelihood,  $\bar{\mathcal{L}}_x^{\max}$  across all  $2K$  models
  - ii. subtract  $\bar{\mathcal{L}}_x^{\max}$  from  $\bar{\mathcal{L}}_{xmk}$  for each model, such that the standardized log likelihood,  $\mathcal{Z}_{xmk}$ , of the best model is 0, and less than 0 for all inferior models.
2. For each model (i.e., each combination of  $m$  and  $k$ ) calculate the mean ( $\bar{\mathcal{Z}}_{mk}$ ) standardized log likelihood of the testing data partition across  $X$  replicates, as well as its standard error ( $SE_{\bar{\mathcal{Z}}_{mk}}$ ) and 95% confidence interval ( $\bar{\mathcal{Z}}_{mk} \pm 1.96 \times SE_{\bar{\mathcal{Z}}_{mk}}$ ).

In other words, the “predictive accuracy” shown as conStruct cross-validation results are in units of improvement in log-likelihood of that model relative to the best model for that partitioning of the data, averaged over data partitions. The standardization is necessary because different data partitions can be systematically more or less difficult to fit, resulting in greater differences in mean training data log likelihood between data partitions than between models fit to the same partition.

If the genomic coordinates of the loci are known, the training/testing partitioning should be designed to accommodate linkage disequilibrium (LD). Loci in strong LD are not inherited independently, so if loci from a single linkage block are included in both training and testing partitions, the independence of the test in the testing partition will be compromised because the parameters estimated from the training partition might be describing process heterogeneity or noise in a region of the genome that also has loci included in the testing partition. The best practice for cross-validation is to make sure that no loci in the testing dataset are in strong LD with, or near on the genome to, loci in the training dataset.

### Calculating layer contributions

Let  $A$  and  $B$  be randomly chosen alleles from samples  $i$  and  $j$  respectively, at a randomly chosen locus. Then, if we let  $U = 2(A - 1/2)$  and  $V = 2(B - 1/2)$ , since  $U$  and  $V$  take the values  $\pm 1$ , so as in Eq. (12),

$$\begin{aligned}\mathbb{E}[UV] &= \mathbb{P}\{U = V\} - \mathbb{P}\{U \neq V\} \\ &= 2\mathbb{P}\{U = V\} - 1 \\ &= 2\mathbb{P}\{A = B\} - 1\end{aligned}$$

To translate,  $\mathbb{P}\{U = V\}$  is the probability that the alleles from our two focal samples agree with each other, while  $\mathbb{P}\{U \neq V\}$  is the probability that they disagree. This implies that  $\mathbb{E}[UV] = 1 - 2\pi_{ij}$ , where  $\pi_{ij}$  is the probability that two randomly chosen alleles differ, which is the genetic divergence.

Now, here is a generative model that gives us the form of the covariance we have postulated. To decide whether or not  $A$  and  $B$  will agree, first each sample randomly chooses a layer: call these layers  $I$  and  $J$ . The probability that  $A$  chooses layer  $k$  is  $\mathbb{P}(I = k) = w_i^{(k)}$ , the  $i$ th sample's admixture proportion in the  $k$ th layer. The same holds true for  $B$ . If they do not choose the same layer, the probability that they agree is  $p_\gamma$ . If they do choose the same layer, then they agree with a probability  $1/2 + p_\gamma + q_{ij}^{(k)}$  that depends on their distance apart. By the above, the probability of agreement is  $\mathbb{P}\{A = B\} = 2 \text{cov}[A, B] + 1/2$ , and so we can define

$$\begin{aligned}p_\gamma &= 2(\gamma + \delta_{ij}\eta_i) \\ q_{ij}^{(k)} &= 2\alpha_0^{(k)} \exp\left(-\left(\alpha_D^{(k)} D_{ij}\right)^{\alpha_2^{(k)}}\right) + 2\phi^{(k)}.\end{aligned}$$

One way to summarize the contribution of each layer is to partition the probability of agreement into contributions due to agreement "in" each layer. So, the contribution from layer  $k$  to agreement between  $i$  and  $j$  is

$$\frac{w_i^{(k)} w_j^{(k)} (1/2 + q_{ij}^{(k)})}{\left(1 - \sum_{k=1}^K w_i^{(k)} w_j^{(k)}\right) (1/2 + p_\gamma) + \sum_{k=1}^K w_i^{(k)} w_j^{(k)} (1/2 + q_{ij}^{(k)})},$$

which is the probability, given that they agree, that they agree thanks to layer  $k$ . Because our signal comes from *variation* in covariance, we omit the  $p_\gamma$  terms (i.e., we condition on agreement not due to "background" levels of agreement in the interpretations above). Stated in this way, this quantity is the relative contribution of the  $k^{\text{th}}$  layer to the (model-based) kinship coefficient between  $i$  and  $j$ .

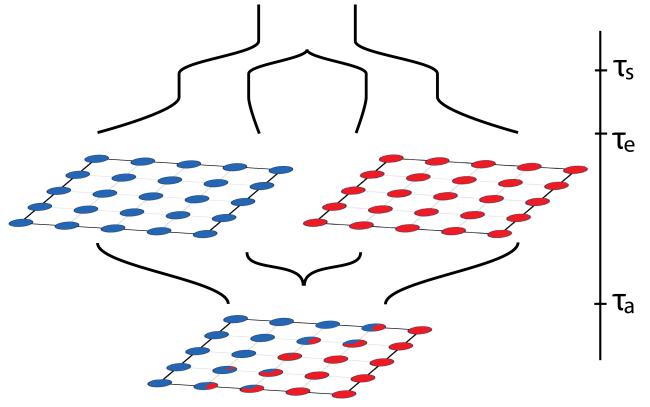
This suggests defining the overall contribution of layer  $k$  to agreement,  $\xi^{(k)}$ , to be the average of that quantity over  $i$  and  $j$ :

$$\xi^{(k)} = \sum_{i=1}^N \sum_{j=i}^N w_i^{(k)} w_j^{(k)} \left( 2\alpha_0^{(k)} \left( \exp\left(-(\alpha_D^{(k)} D_{i,j})^{\alpha_2^{(k)}}\right) \right) + 2\phi^{(k)} + \frac{1}{2} \right), \quad (16)$$

which is that layer's contribution to agreement between samples summed over the upper triangle (excluding the diagonal) of the covariance matrix. We define the contribution of the  $k$ th layer,  $\Xi^{(k)}$ , as the relative contribution of the  $k$ th layer to total agreement:

$$\Xi^{(k)} = \frac{\xi^{(k)}}{\sum_{k=1}^K \xi^{(k)}}. \quad (17)$$

This is the quantity that is plotted in Figures 4 and 9.



**Figure 10** Schematic of how we simulate datasets with continuous and discrete differentiation, using  $K = 2$  as an example. Going forward in time, the  $K$  populations split from a common ancestor at time  $\tau_s$ , then expand to each colonize a lattice of demes with nearest-neighbor symmetric migration at time  $\tau_e$ , then finally at time  $\tau_a$  collapse into a single lattice consisting of demes with ancestry entirely in one or the other of the populations, or admixed between them.

### Simulation details

We wished to simulate data under a model that had some biological realism, but at the same time had unambiguous true admixture proportions (so as to test the behavior of the method). This second requirement precluded scenarios of, e.g. recent secondary contact between populations expanding out of different refugia, which would have more biological realism, but no unambiguous ancestry proportions for admixed populations. Here, we describe in more detail the procedure we use to simulate our test dataset, using a cartoon schematic with  $K = 2$  as an example (Fig 10).

Using the program `ms` (Hudson 2002), we generated discrete population structure by simulating  $K$  distinct populations, each of which split from a common ancestor  $\tau_s$  units of coalescent time in the past, without subsequent migration between them. Then, to generate continuous differentiation within each population, at time  $\tau_e$  in the past, each of these discrete populations instantaneously colonizes an independent lattice of demes, for which we use a stepping stone model with symmetric migration to nearest neighbors (eight neighbors, including diagonals).

Finally, at time  $\tau_a$  in the past we generate a single dataset by collapsing those  $K$  discrete lattices into a single grid of demes that are admixed to various degrees from these different layers. We wish to simulate realistic patterns of admixture (and thereby set a more difficult test for the method), by generating spatially autocorrelated admixture proportions in each diverged population. To do so, we first place  $K$  equidistant points on the circle centered on our lattice. These points serve as "foci" of ancestry in each of the  $K$  layers. We then calculate the distance from each deme in the sampled lattice to each of these  $K$  foci, and draw admixture proportions for each deme from a Dirichlet distribution for which the concentration parameter for deme  $i$  in layer  $k$  is inversely proportional to the distance between deme  $i$  and focus  $k$ . This creates a pattern in which the admixture proportions in a given layer decreases with the distance from that layer's focus, as might be expected if a spatial process were mediating admixture between diverged populations.

The parameters used to simulate the data were as follows: a diploid population size of 1000, a migration rate between neighboring demes of 0.4, a deep split time between layers of 500 (corresponding to  $\tau_s$  in Fig 10), an expansion event across layers of 250 (corresponding to  $\tau_e$  in Fig 10), and an admixture event between layers in the immediate past ( $1 \times 10^{-4}$ ). The times and rates reported above have already been scaled by  $4N$  (as per ms syntax), and therefore give the values fed directly to ms. We used the -s option to sample a single segregating site per coalescent history, and simulated  $1 \times 10^4$  independent histories — corresponding to the same number of independent loci — in each dataset, with 10 diploid genotypes generated per deme at each locus.

## Acknowledgments

We thank Marjorie Weber, Yaniv Brandvain, William Wetzel, Mariah Meek, Doc Edge, Evan McCartney-Melstad, Matthew Stephens, Nick Barton, and the anonymous reviewers for invaluable comments on the method and manuscript, as well as Quentin Cronk, who provided input on the *Populus* analyses, and Emily Puckett, who provided input on the black bear analyses. We also thank the attendees at the 2017 SSE Meeting in Portland, OR, whose votes determined the name of the method. This work was supported in part by the National Science Foundation under award number NSF #1262645 (DBI) to PR and GC, the National Institute of General Medical Sciences of the National Institutes of Health under award numbers NIH R01-GM108779 to GC, and the National Science Foundation under award numbers NSF #1148897 and #1402725 to GB.

## Literature Cited

- Alexander, D. H. and K. Lange, 2011 Enhancements to the admixture algorithm for individual ancestry estimation. *BMC Bioinformatics* **12**: 246.
- Alexander, D. H., J. Novembre, and K. Lange, 2009 Fast model-based estimation of ancestry in unrelated individuals. *Genome Research* **19**: 1655–1664.
- Barton, N. H., 2008 The effect of a barrier to gene flow on patterns of geographic variation. *Genetics Research* **90**: 139–149.
- Barton, N. H., F. Depaulis, and A. M. Etheridge, 2002 Neutral evolution in spatially continuous populations. *Theoretical Population Biology* **61**: 31–48.
- Barton, N. H., A. M. Etheridge, and A. Véber, 2013 Modelling evolution in a spatial continuum. *Journal of Statistical Mechanics: Theory and Experiment* **2013**: P01002.
- Bradburd, G. S., P. L. Ralph, and G. M. Coop, 2013 Disentangling the effects of geographic and ecological isolation on genetic differentiation. *Evolution* **67**: 3258–3273.
- Bradburd, G. S., P. L. Ralph, and G. M. Coop, 2016 A spatial framework for understanding population structure and admixture. *PLoS Genet* **12**: 1–38.
- Byun, S. A., B. F. Koop, and T. E. Reimchen, 1997 North american black bear mtDNA phylogeography: Implications for morphology and the haida gwaii glacial refugium controversy. *Evolution* **51**: 1647–1653.
- Carpenter, B., 2015 Stan: A probabilistic programming language. *Journal of Statistical Software*.
- Cavalli-Sforza, L. L. and A. Piazza, 1975 Analysis of evolution: Evolutionary rates, independence and treeness. *Theoretical Population Biology* **8**: 127 – 165.
- Caye, K., F. Jay, O. Michel, and O. François, 2018 Fast inference of individual admixture coefficients using geographic data. *The Annals of Applied Statistics* **12**: 586–608.
- Chang-Yi, X., C. M. R., and Y. C. C., 2012 Ecotypic mode of regional differentiation of black cottonwood (*populus trichocarpa*) due to restricted gene migration: further evidence from a field test on the northern coast of british columbia. *Canadian Journal of Forest Research* **42**: 400–405.
- Corander, J., P. Waldmann, and M. J. Sillanpää, 2003 Bayesian analysis of genetic differentiation between populations. *Genetics* **163**: 367–374.
- Cronk, Q. C. B., 2005 Plant eco-devo: the potential of poplar as a model organism. *New Phytologist* **166**: 39–48.
- Diggle, P. J., J. A. Tawn, and R. A. Moyeed, 1998 Model-based geostatistics. *Jounal of the Royal Statistical Society. Series C (Applied Statistics)* **47**: 299–350.
- Eckenwalder, J. E., 1984 Natural intersectional hybridization between north american species of *populus* (salicaceae) in sections aigeiros and tacamahaca. ii. taxonomy. *Canadian Journal of Botany* **62**: 325–335.
- Engelhardt, B. E. and M. Stephens, 2010 Analysis of population structure: A unifying framework and novel methods based on sparse factor analysis. *PLOS Genetics* **6**: 1–12.
- Epperson, B., 2003 *Geographical Genetics*. Monographs in Population Biology, Princeton University Press.
- Evanno, G., S. Regnaut, and J. Goudet, 2005 Detecting the number of clusters of individuals using the software structure: a simulation study. *Molecular Ecology* **14**: 2611–2620.
- Falush, D., M. Stephens, and J. K. Pritchard, 2003 Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. *Genetics* **164**: 1567–1587.
- Falush, D., L. van Dorp, and D. Lawson, 2016 A tutorial on how (not) to over-interpret STRUCTURE/ADMIIXTURE bar plots. *bioRxiv* .
- Frantz, A. C., S. Cellina, A. Krier, L. Schley, and T. Burke, 2009 Using spatial bayesian methods to determine the genetic structure of a continuously distributed population: clusters or isolation by distance? *Journal of Applied Ecology* **46**: 493–505.
- Geraldes, A., S. P. DiFazio, G. T. Slavov, P. Ranjan, W. Muchero, et al., 2013 A 34k SNP genotyping array for *populus trichocarpa*: Design, application to the study of natural populations and transferability to other *populus* species. *Molecular Ecology Resources* **13**: 306–323.
- Geraldes, A., N. Farzaneh, C. J. Grassa, A. D. McKown, R. D. Guy, et al., 2014 Landscape genomics of *populus trichocarpa*: the role of hybridization, limited gene flow, and natural selection in shaping patterns of population structure. *Evolution* **68**: 3260–3280.
- Guillot, G., F. Mortier, and A. Estoup, 2005 Geneland: a computer package for landscape genetics. *Molecular Ecology Notes* **5**: 712–715.
- Haak, W., I. Lazaridis, N. Patterson, N. Rohland, S. Mallick, et al., 2015 Massive migration from the steppe was a source for indo-european languages in europe. *Nature* **522**: 207–211.
- Hoffman, M. D. and A. Gelman, 2014 The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* .
- Hubisz, M. J., D. Falush, M. Stephens, and J. K. Pritchard, 2009 Inferring weak population structure with the assistance of sample group information. *Molecular Ecology Resources* **9** *5*: 1322–32.
- Hudson, R. R., 2002 Generating samples under a Wright-Fisher

- neutral model of genetic variation. *Bioinformatics* **18**: 337–338.
- Huelsenbeck, J. P. and P. Andolfatto, 2007 Inference of population structure under a Dirichlet process model. *Genetics* **175**: 1787–1802.
- Irwin, D. E., S. Bensch, and T. D. Price, 2001 Speciation in a ring. *Nature* **409**: 333–337.
- Joseph, T. A. and I. Pe'er, 2018 Inference of population structure from ancient DNA. *bioRxiv*.
- Keller, S. R., M. S. Olson, S. Slim, W. Schroeder, and P. Tiffin, 2010 Genomic diversity, population structure, and migration following rapid range expansion in the balsam poplar, *Populus balsamifera*. *Molecular Ecology* **19**: 1212–1226.
- Kimura, M. and G. H. Weiss, 1964 The stepping stone model of population structure and the decrease of genetic correlation with distance. *Genetics* **49**: 561–576.
- Lawson, D. J., G. Hellenthal, S. Myers, and D. Falush, 2012 Inference of population structure using dense haplotype data. *PLoS Genet* **8**: e1002453.
- Lazaridis, I., N. Patterson, A. Mittnik, G. Renaud, S. Mallick, et al., 2014 Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* **513**: 409–413.
- Linck, E. B. and C. J. Battey, 2017 Minor allele frequency thresholds strongly affect population structure inference with genomic datasets. *bioRxiv*.
- Malécot, G., 1969 *The Mathematics of Heredity*. Freeman, Translated from the French edition, 1948.
- McKown, A. D., R. D. Guy, J. Klapste, A. Geraldes, M. Friedmann, et al., 2014 Geographical and environmental gradients shape phenotypic trait variation and genetic structure in *Populus trichocarpa*. *New Phytologist* **201**: 1263–1276.
- McVean, G., 2009 A genealogical interpretation of principal components analysis. *PLoS Genet* **5**: e1000686.
- Meirmans, P., 2009 Genodive version 2.0 b14. Computer software distributed by the author. Available from: <http://www.bentleydrummer.nl/software/software/GenoDive.html>.
- Meirmans, P. G., 2012 The trouble with isolation by distance. *Molecular Ecology* **21**: 2839–2846.
- Menozzi, P., A. Piazza, and L. Cavalli-Sforza, 1978 Synthetic maps of human gene frequencies in Europeans. *Science* **201**: 786–792.
- Moritz, C., 1994 Defining “evolutionarily significant units” for conservation. *Trends in Ecology and Evolution* **9**: 373 – 375.
- Moritz, C., V. Funk, and A. K. Sakai, 2002 Strategies to protect biological diversity and the evolutionary processes that sustain it. *Systematic Biology* **51**: 238–254.
- Nagylaki, T., 1978 A diffusion model for geographically structured populations. *Journal of Mathematical Biology* **6**: 375–382.
- Nagylaki, T. and V. Barcilon, 1988 The influence of spatial inhomogeneities on neutral models of geographical variation: II. the semi-infinite linear habitat. *Theoretical Population Biology* **33**: 311 – 343.
- Nielsen, R., J. M. Akey, M. Jakobsson, J. K. Pritchard, S. Tishkoff, et al., 2017 Tracing the peopling of the world through genomics. *Nature* **541**: 302–310.
- Novembre, J. and M. Stephens, 2008 Interpreting principal component analyses of spatial population genetic variation. *Nature Genetics* **40**: 646–649.
- Patterson, N., P. Moorjani, Y. Luo, S. Mallick, N. Rohland, et al., 2012 Ancient admixture in human history. *Genetics* **192**: 1065–1093.
- Patterson, N., A. L. Price, and D. Reich, 2006 Population structure and eigenanalysis. *PLoS Genet* **2**: e190.
- Peter, B. M., 2016 Admixture, population structure and f-statistics. *Genetics* .
- Petkova, D., J. Novembre, and M. Stephens, 2016 Visualizing spatial population structure with estimated effective migration surfaces. *Nature Genetics* **48**: 94–100C.
- Picard, R. R. and R. D. Cook, 1984 Cross-validation of regression models. *Journal of the American Statistical Association* **79**: 575–583.
- Pickrell, J. K. and J. K. Pritchard, 2012 Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genet* **8**: e1002967.
- 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. *Nat Genet* **38**: 904–909.
- Pritchard, J. K., M. Stephens, and P. Donnelly, 2000 Inference of population structure using multilocus genotype data. *Genetics* **155**: 945–959.
- Puckett, E. E., P. D. Etter, E. A. Johnson, and L. S. Eggert, 2015 Phylogeographic analyses of american black bears (*Ursus americanus*) suggest four glacial refugia and complex patterns of postglacial admixture. *Molecular Biology and Evolution* **32**: 2338–2350.
- Raj, A., M. Stephens, and J. K. Pritchard, 2014 fastSTRUCTURE: Variational inference of population structure in large SNP data sets. *Genetics* **197**: 573–589.
- Ringbauer, H., A. Kolesnikov, D. L. Field, and N. H. Barton, 2018 Estimating barriers to gene flow from distorted isolation-by-distance patterns. *Genetics* **208**: 1231–1245.
- Rosenberg, N. A., S. Mahajan, S. Ramachandran, C. Zhao, J. K. Pritchard, et al., 2005 Clines, clusters, and the effect of study design on the inference of human population structure. *PLoS Genet* **1**.
- Sawyer, S., 1976 Results for the stepping stone model for migration in population genetics. *The Annals of Probability* **4**: 699–728.
- Schraiber, J., 2017 Assessing the relationship of ancient and modern populations. *bioRxiv* .
- Serre, D. and S. Pääbo, 2004 Evidence for gradients of human genetic diversity within and among continents. *Genome Research* **14**: 1679–1685.
- Sexton, J. P., S. B. Hangartner, and A. A. Hoffmann, 2014 Genetic isolation by environment or distance: which pattern of gene flow is most common? *Evolution* **68**: 1–15.
- Shiga, T., 1988 Stepping stone models in population genetics and population dynamics. In *Stochastic processes in physics and engineering (Bielefeld, 1986)*, volume 42 of *Math. Appl.*, pp. 345–355, Reidel, Dordrecht.
- Slatkin, M., 1985 Gene flow in natural populations. *Annual Review of Ecology and Systematics* **16**: 393–430.
- Slatkin, M. and F. Racimo, 2016 Ancient dna and human history. *Proceedings of the National Academy of Sciences* **113**: 6380–6387.
- Slavov, G. T., S. P. DiFazio, J. Martin, W. Schackwitz, W. Muchero, et al., 2012 Genome resequencing reveals multiscale geographic structure and extensive linkage disequilibrium in the forest tree *Populus trichocarpa*. *New Phytologist* **196**: 713–725.
- Stan Development Team, 2015 Stan: A C++ library for probability and sampling, version 2.10.0.
- Stan Development Team, 2016 RStan: the R interface to Stan, version 2.10.1.

- Stone, K. D. and J. A. Cook, 2000 Phylogeography of black bears (*ursus americanus*) of the pacific northwest. Canadian Journal of Zoology **78**: 1218–1223.
- Suarez-Gonzalez, A., C. A. Hefer, C. Christe, O. Corea, C. Lexer, *et al.*, 2016 Genomic and functional approaches reveal a case of adaptive introgression from *populus balsamifera* (balsam poplar) in *p. trichocarpa* (black cottonwood). Molecular Ecology **25**: 2427–2442.
- Verity, R. and R. A. Nichols, 2016 Estimating the number of subpopulations (K) in structured populations. Genetics **203**: 1827–1839.
- Wake, D. B. and C. J. Schneider, 1998 Taxonomy of the plethodontid salamander genus *ensatina*. Herpetologica **54**: pp. 279–298.
- Waples, R., 1998 Separating the wheat from the chaff: patterns of genetic differentiation in high gene flow species. Journal of Heredity **89**: 438–450.
- Wasser, S. K., A. M. Shedlock, K. Comstock, E. Ostrander, B. Mutayoba, *et al.*, 2004 Assigning African elephant DNA to geographic region of origin: Applications to the ivory trade. PNAS **101**: 14847–52.
- Wooding, S. and R. Ward, 1997 Phylogeography and pleistocene evolution in the north american black bear. Molecular Biology and Evolution **14**: 1096–1105.
- Wright, S., 1943 Isolation by distance. Genetics **28**: 114–138.
- Wright, S., 1949 The genetical structure of populations. Annals of Eugenics **15**: 323–354.
- Xie, C.-Y., C. C. Ying, A. D. Yanchuk, and D. L. Holowachuk, 2009 Ecotypic mode of regional differentiation caused by restricted gene migration: a case in black cottonwood (*populus trichocarpa*) along the pacific northwest coast. Canadian Journal of Forest Research **39**: 519–525.
- Zhou, L. and J. A. Holliday, 2012 Targeted enrichment of the black cottonwood (*populus trichocarpa*) gene space using sequence capture. BMC Genomics **13**: 703.

## **Supplementary Materials**

# **Supplementary Figures for: Inferring continuous and discrete population genetic structure across space**

Gideon S. Bradburd<sup>1,\*</sup>, Graham M. Coop<sup>2,❸</sup>, Peter L. Ralph<sup>3,❸</sup>,

**1** Ecology, Evolutionary Biology, and Behavior Graduate Group Department of Integrative Biology, Michigan State University, MI 48824

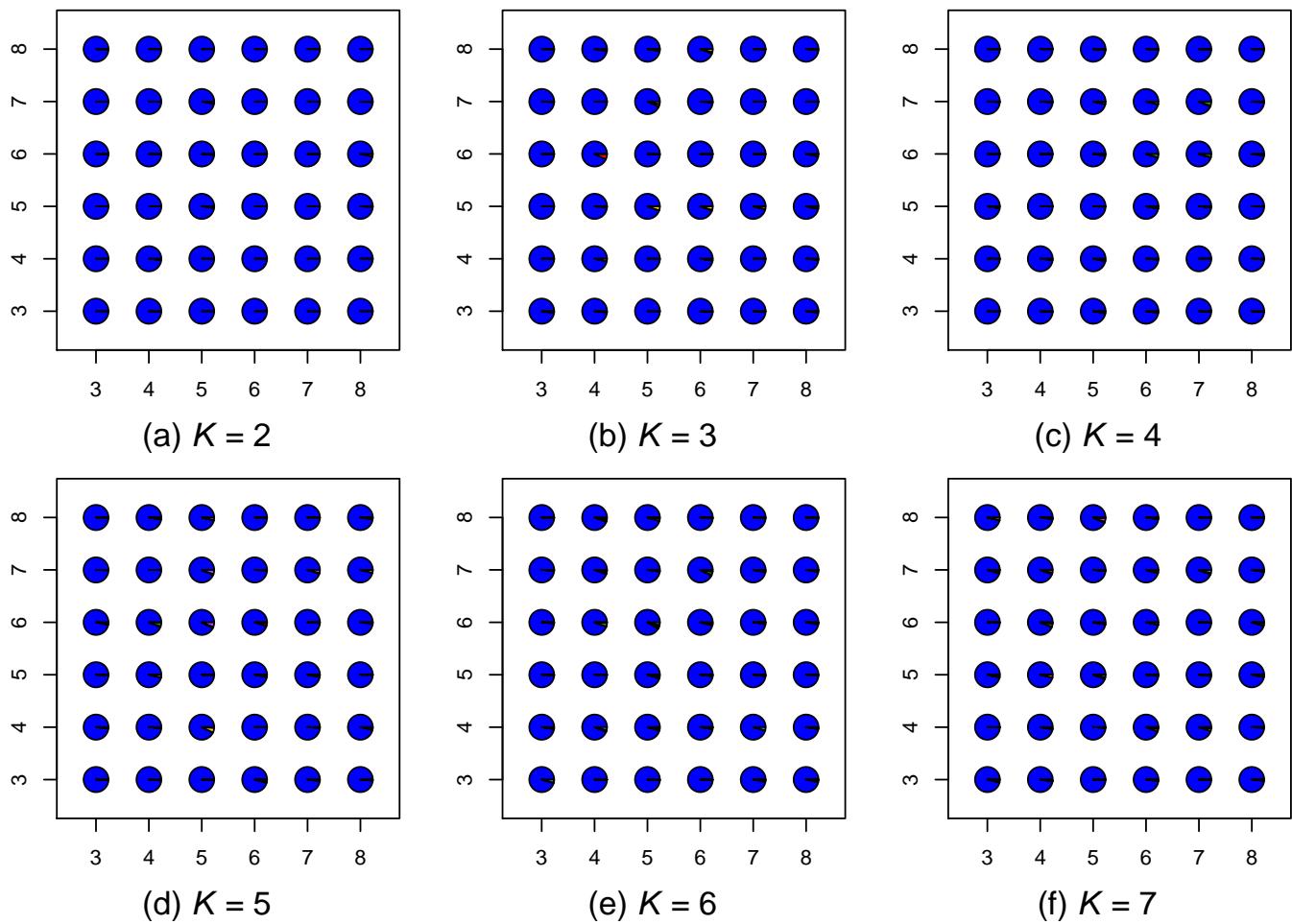
**2** Center for Population Biology, Department of Evolution and Ecology, University of California, Davis, CA 95616

**3** Institute of Ecology and Evolution, Departments of Mathematics and Biology, University of Oregon, Eugene, OR 97403

\*bradburd@msu.edu

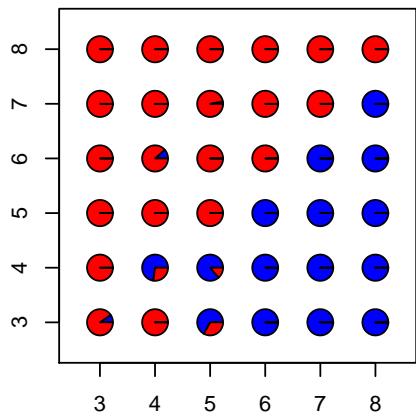
❸These authors contributed equally to this work.

## True $K = 1$

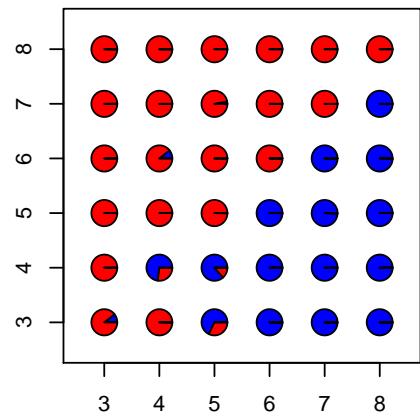


**Figure S1** Map of admixture proportions estimated using a spatial conStruct model for  $K = 2$  through 7. The data were simulated using one layer with nearest-neighbor symmetric migration between demes.

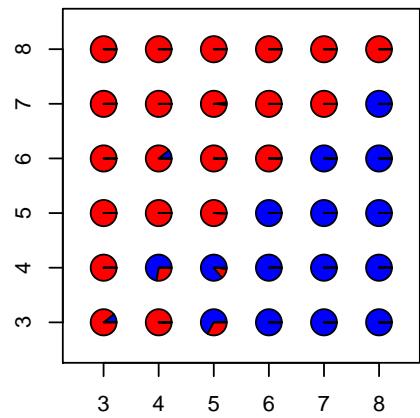
## True $K = 2$



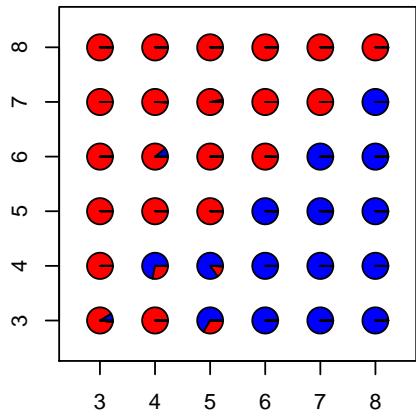
(a)  $K = 2$



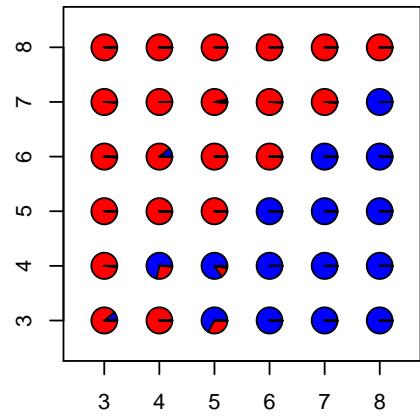
(b)  $K = 3$



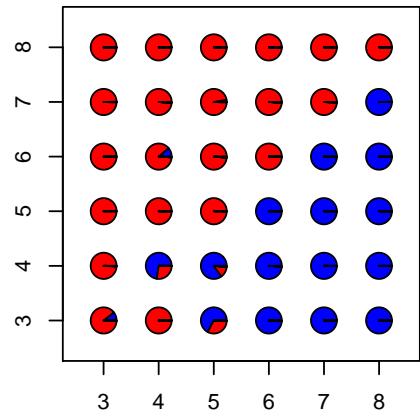
(c)  $K = 4$



(d)  $K = 5$



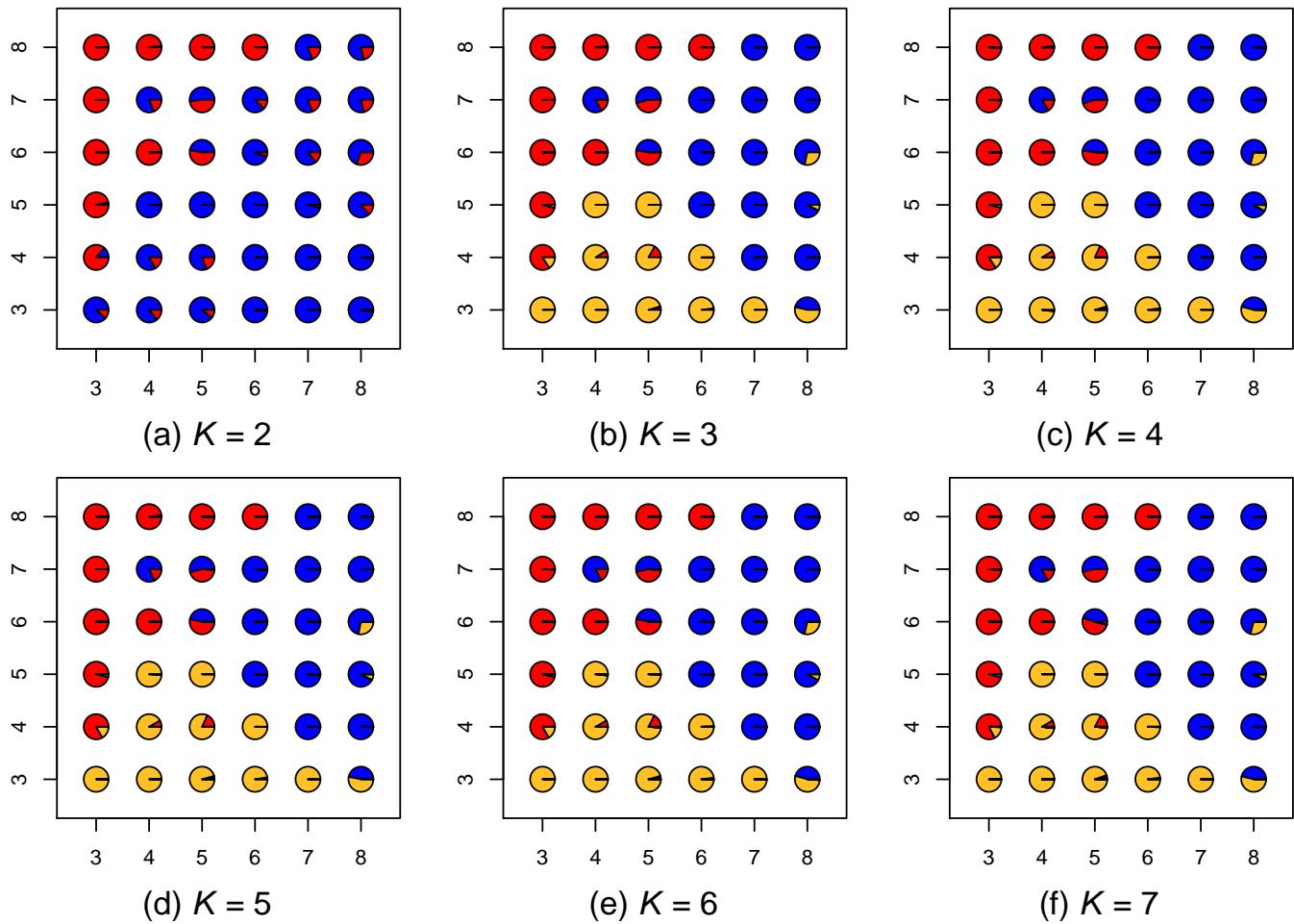
(e)  $K = 6$



(f)  $K = 7$

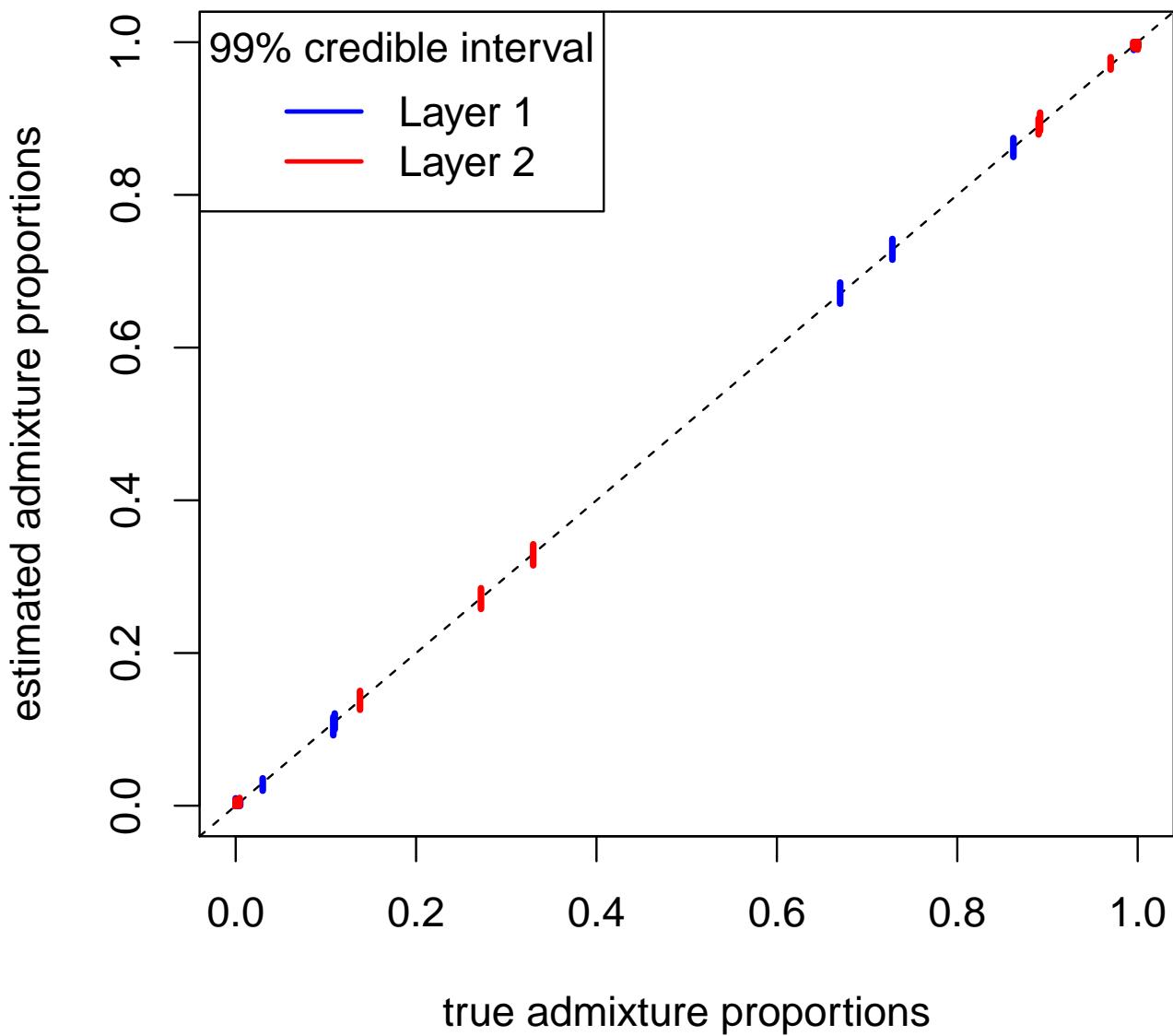
**Figure S2** Map of admixture proportions estimated using a spatial conStruct model for  $K = 2$  through 7. The data were simulated using two layers with nearest-neighbor symmetric migration between demes.

### True $K = 3$



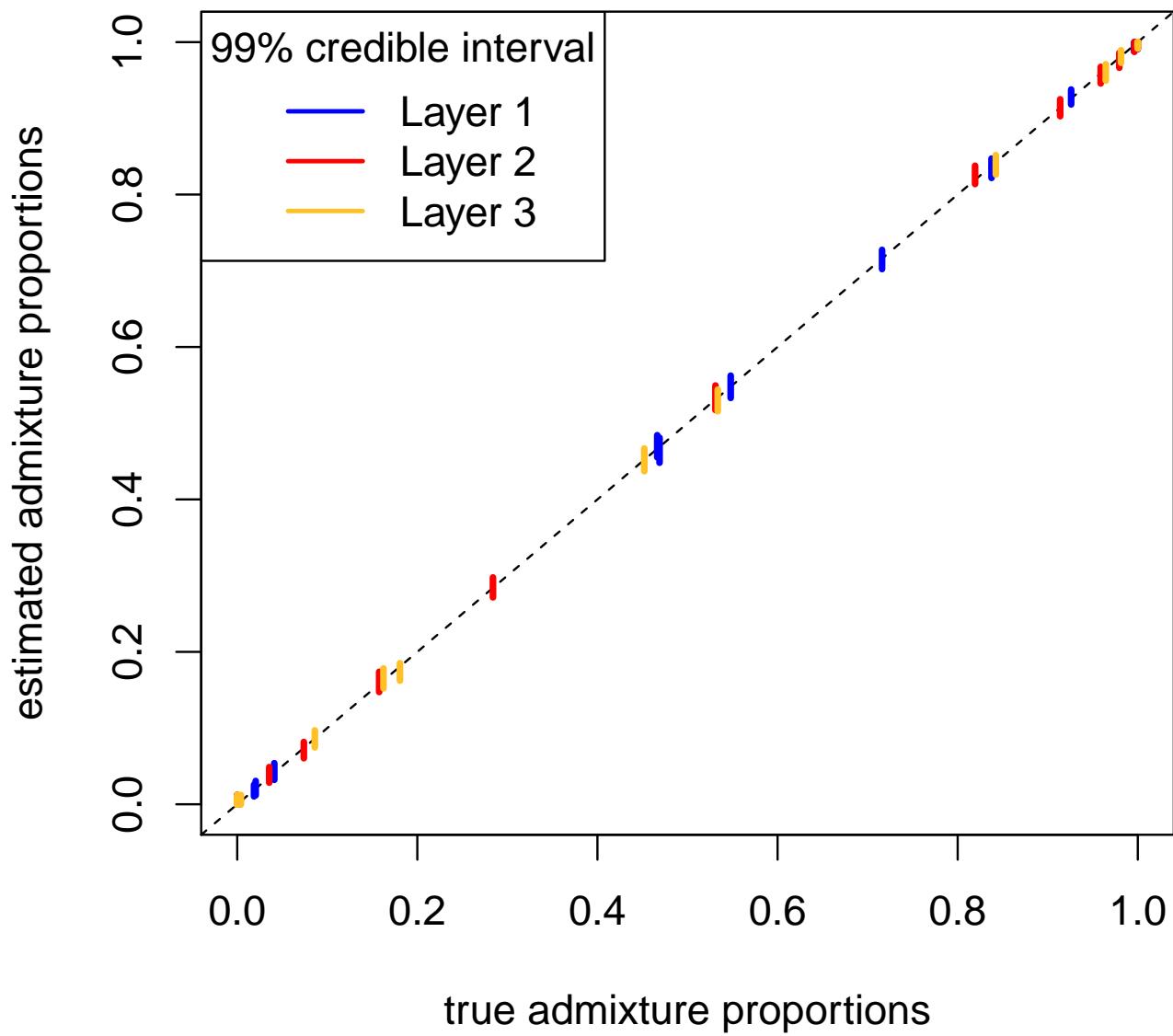
**Figure S3** Map of admixture proportions estimated using a spatial conStruct model for  $K = 2$  through 7. The data were simulated using three layers with nearest-neighbor symmetric migration between demes.

## Fitting admixture parameters (true $K = 2$ )



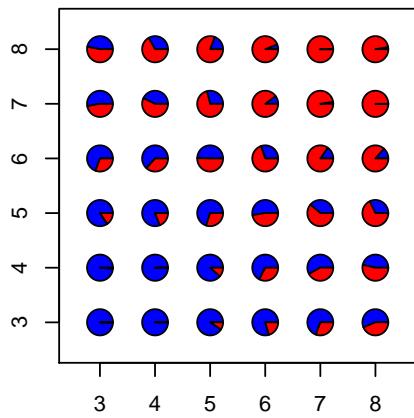
**Figure S4** Plot of conStruct's ability to correctly estimate admixture proportions on simulated data, from an analysis with a spatial model using  $K = 2$ . The horizontal axis shows the admixture proportions used to simulate the data, and the vertical axis shows the 99% credible intervals for those proportions as reported by conStruct.

## Fitting admixture parameters (true $K = 3$ )

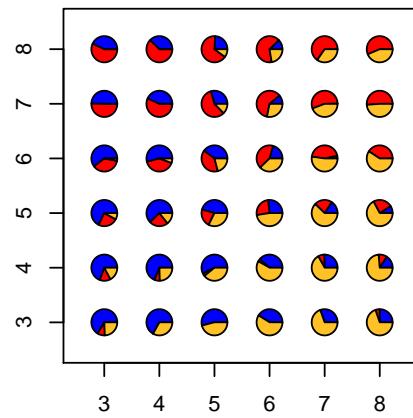


**Figure S5** Plot of conStruct ability to correctly estimate admixture proportions on simulated data, from an analysis with a spatial model using  $K = 3$ . The horizontal axis shows the admixture proportions used to simulate the data, and the vertical axis shows the 99% credible intervals for those proportions as reported by conStruct.

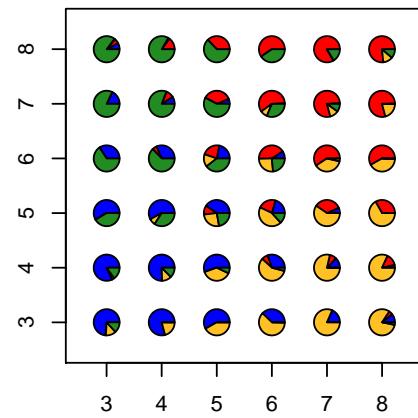
## True $K = 1$



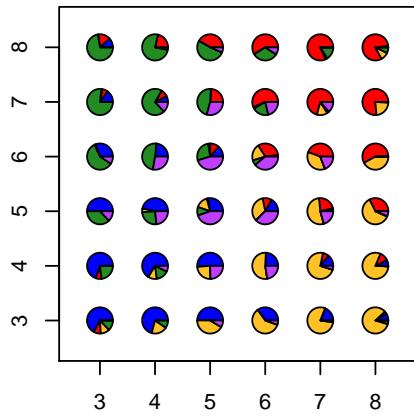
(a)  $K = 2$



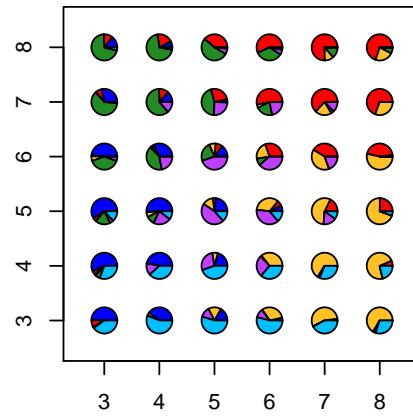
(b)  $K = 3$



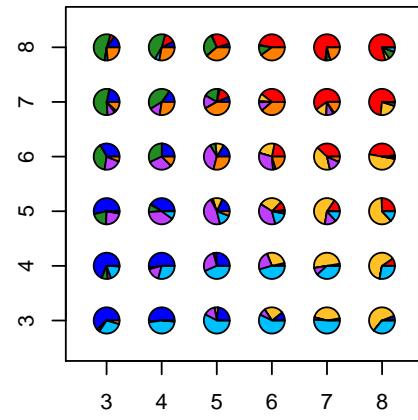
(c)  $K = 4$



(d)  $K = 5$



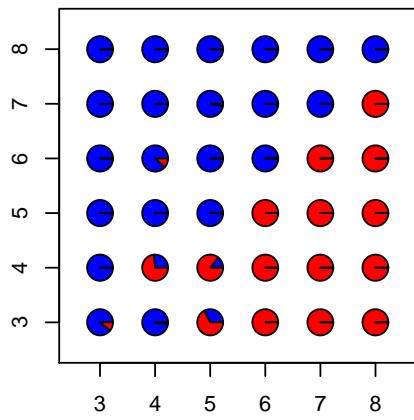
(e)  $K = 6$



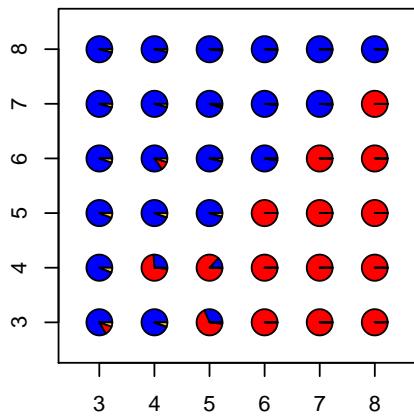
(f)  $K = 7$

**Figure S6** Map of admixture proportions estimated using a nonspatial conStruct model for  $K = 2$  through 7. The data were simulated using one layer with nearest-neighbor symmetric migration between demes.

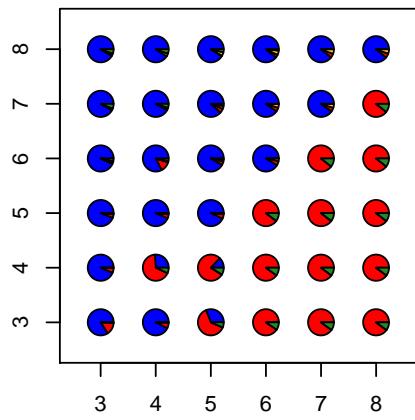
## True $K = 2$



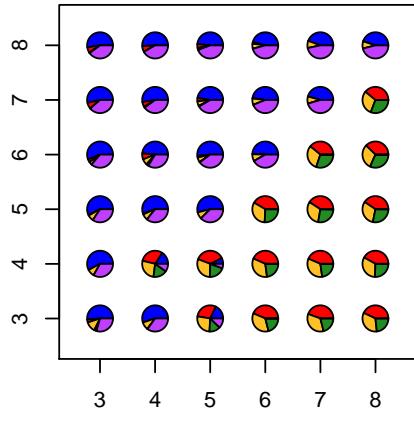
(a)  $K = 2$



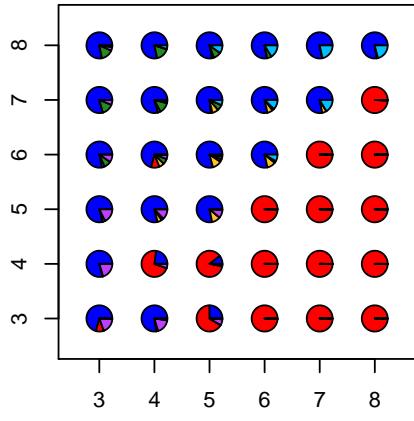
(b)  $K = 3$



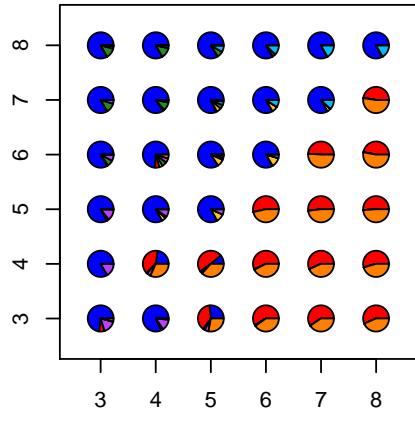
(c)  $K = 4$



(d)  $K = 5$



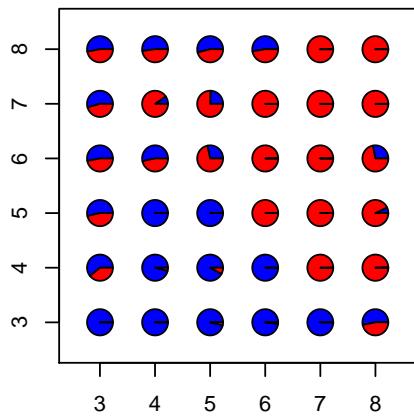
(e)  $K = 6$



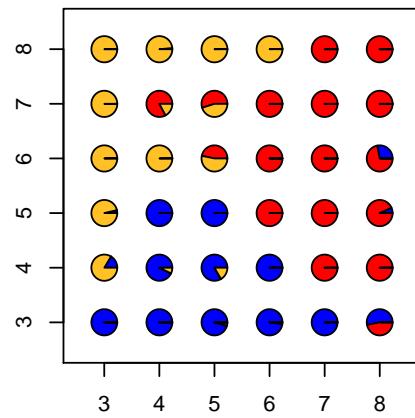
(f)  $K = 7$

**Figure S7** Map of admixture proportions estimated using a nonspatial conStruct model for  $K = 2$  through 7. The data were simulated using two layers with nearest-neighbor symmetric migration between demes.

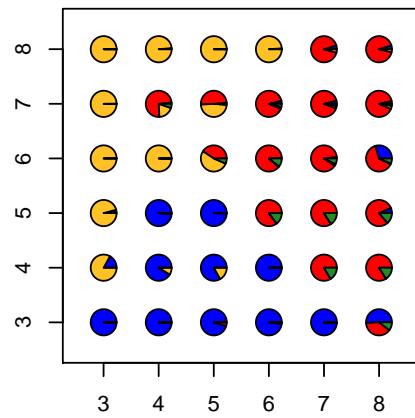
## True $K = 3$



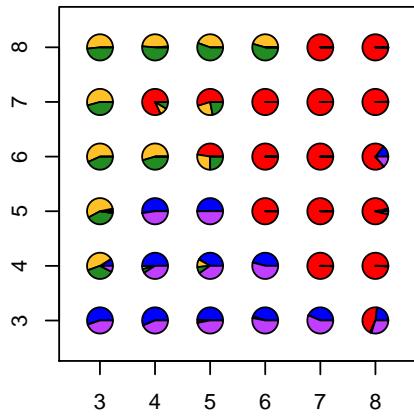
(a)  $K = 2$



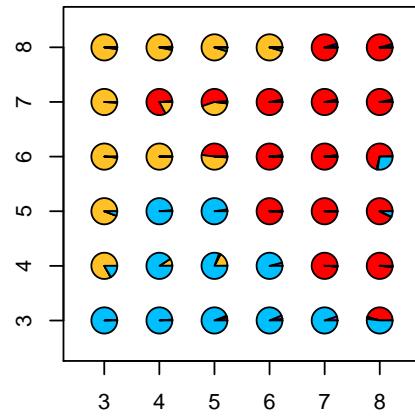
(b)  $K = 3$



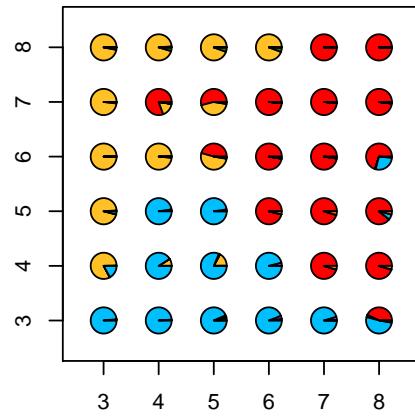
(c)  $K = 4$



(d)  $K = 5$

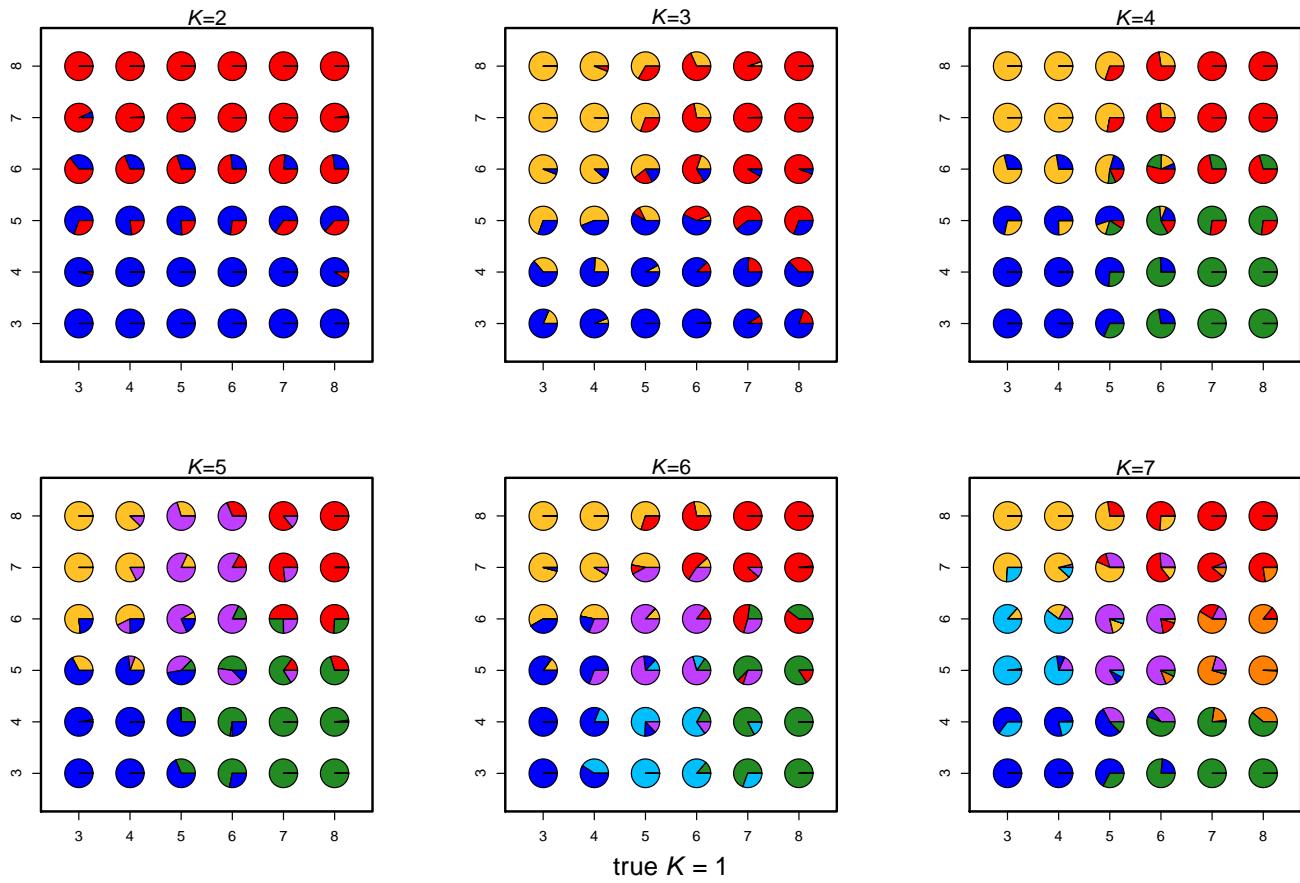


(e)  $K = 6$

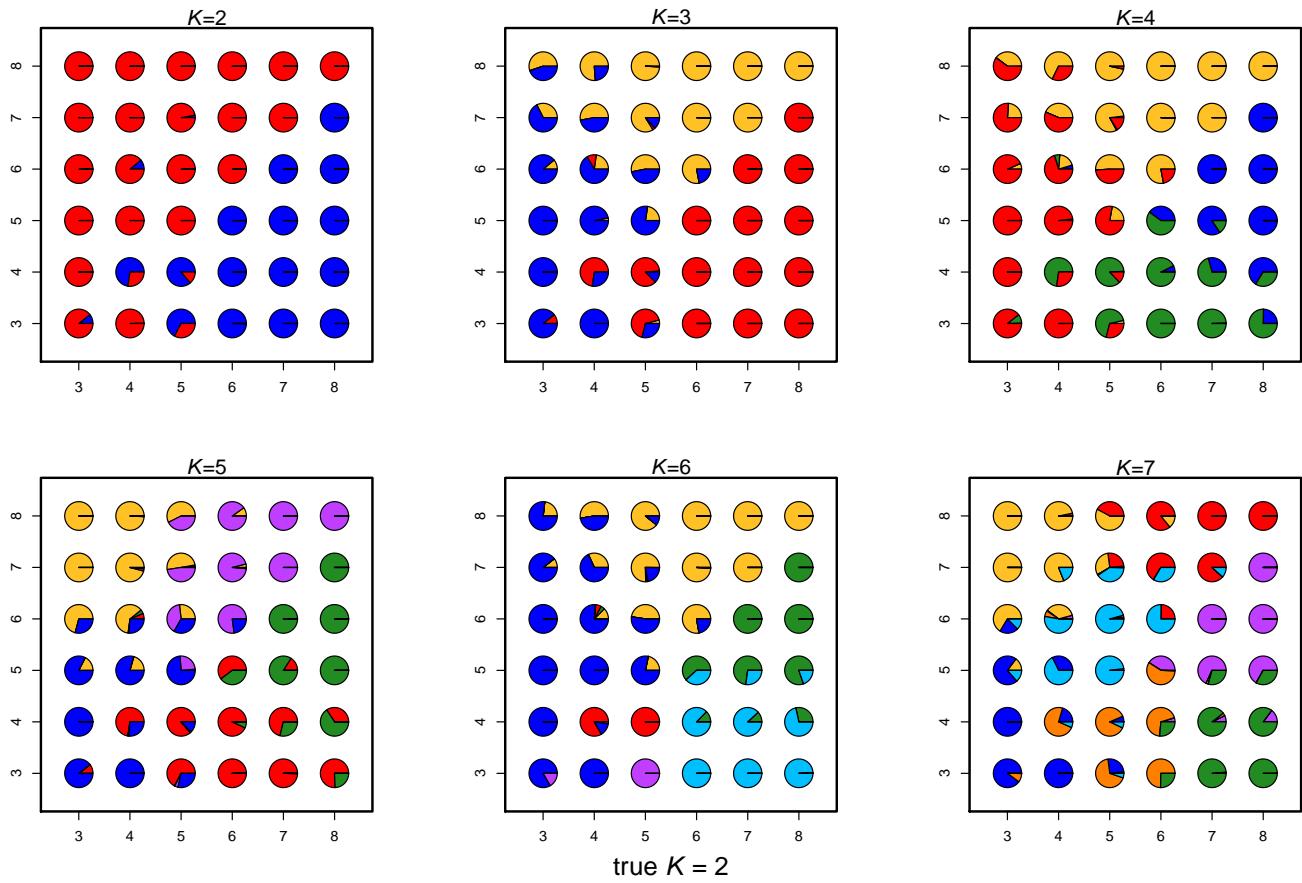


(f)  $K = 7$

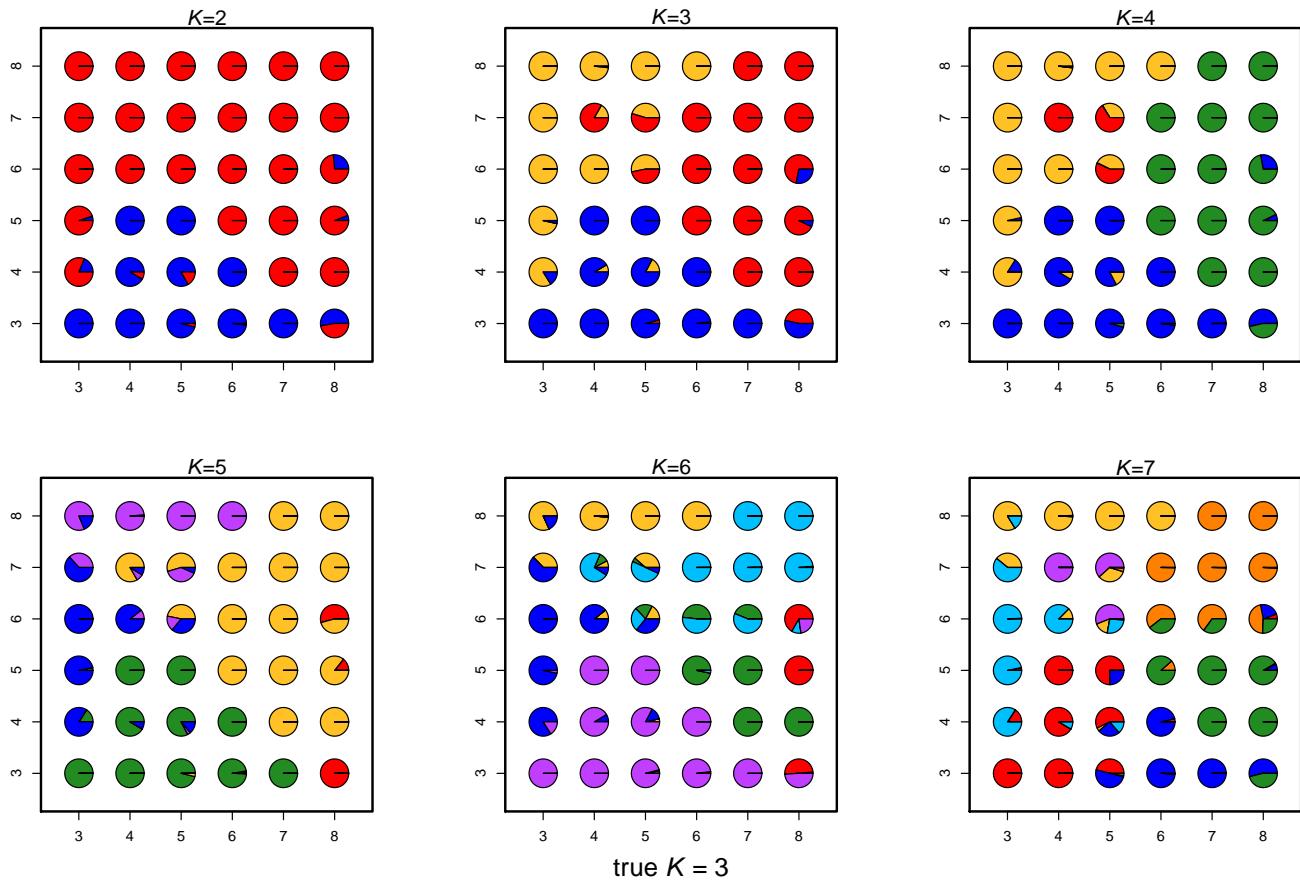
**Figure S8** Map of admixture proportions estimated using a nonspatial conStruct model for  $K = 2$  through 7. The data were simulated using three layers with nearest-neighbor symmetric migration between demes.



**Figure S9** Map of admixture proportions estimated using ADMIXTURE [Alexander et al. \(2009\)](#) for  $K = 2$  through 7. The data were simulated using one layer with nearest-neighbor symmetric migration between demes. The true value was  $K = 1$ , but the model with the lowest cross-validation error (i.e., the preferred model) was  $K = 7$ .

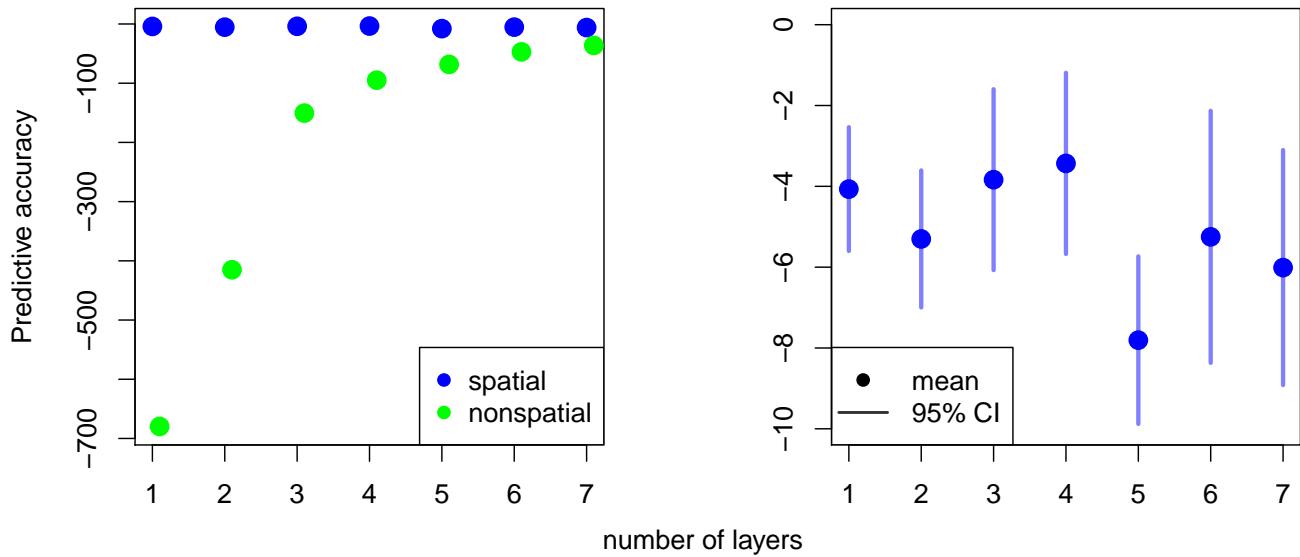


**Figure S10** Map of admixture proportions estimated using ADMIXTURE [Alexander et al. \(2009\)](#) for  $K = 2$  through 7. The data were simulated using two layers with nearest-neighbor symmetric migration between demes. The true value was  $K = 2$ , but the model with the lowest cross-validation error (i.e., the preferred model) was  $K = 7$ .



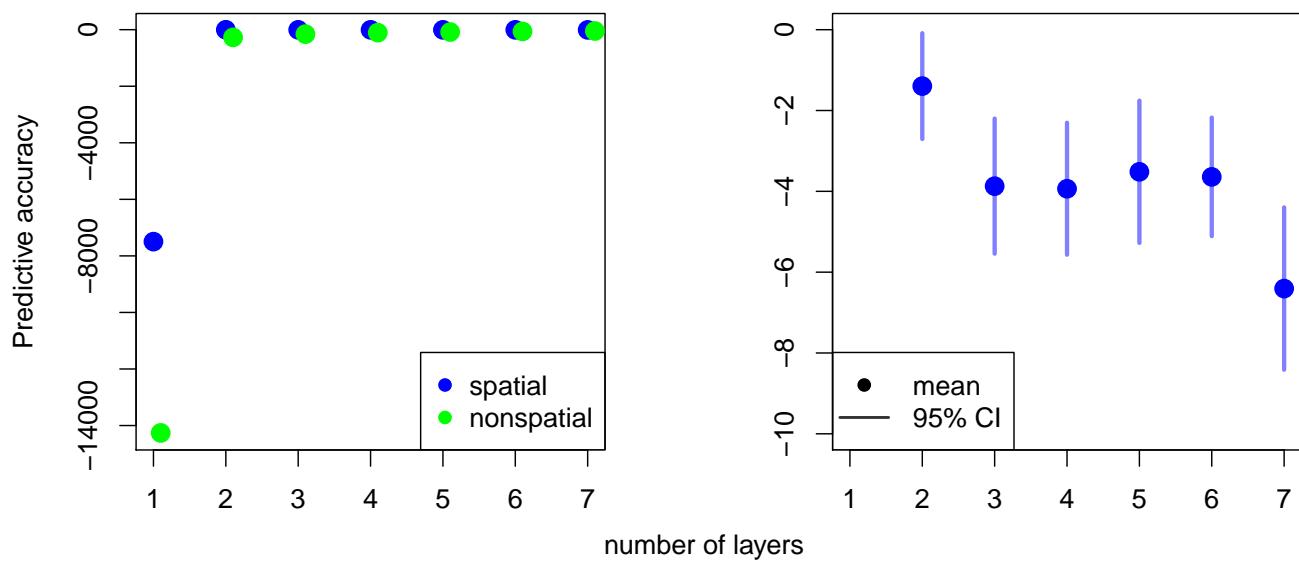
**Figure S11** Map of admixture proportions estimated using ADMIXTURE [Alexander et al. \(2009\)](#) for  $K = 2$  through 7. The data were simulated using three layers with nearest-neighbor symmetric migration between demes. The true value was  $K = 3$ , but the model with the lowest cross-validation error (i.e., the preferred model) was  $K = 7$ .

### Cross-validation results (true $K=1$ )



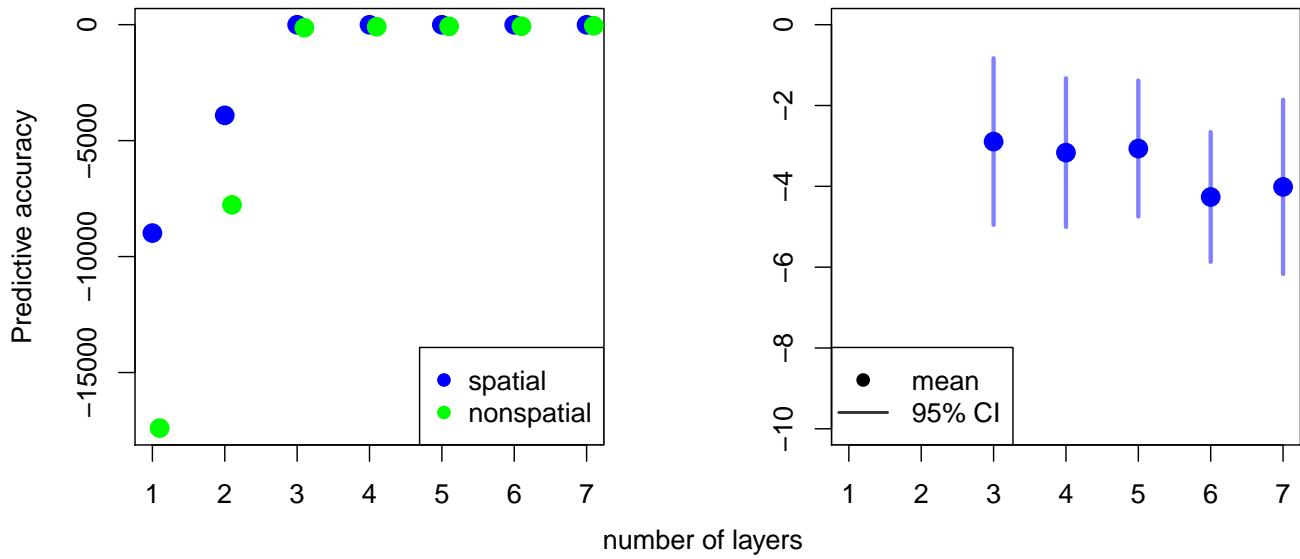
**Figure S12** Cross-validation results for data simulated under  $K = 1$ , comparing the spatial and nonspatial `conStruct` models run with  $K = 1$  through 7. The right panel zooms in on just the spatial cross-validation results.

### Cross-validation results (true $K=2$ )



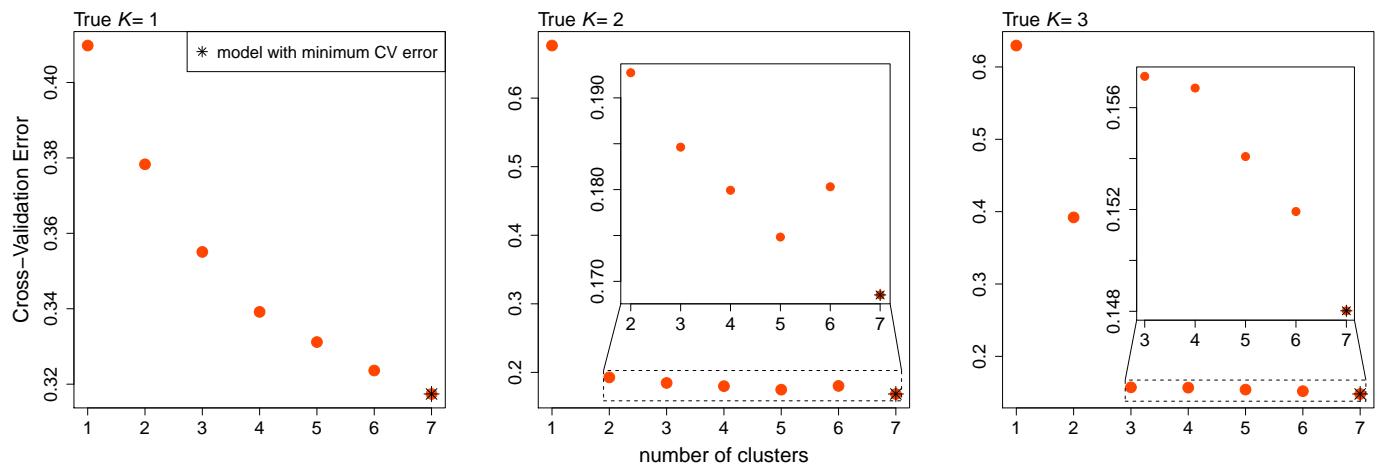
**Figure S13** Cross-validation results for data simulated under  $K = 2$ , comparing the spatial and nonspatial `conStruct` models run with  $K = 1$  through 7. The right panel zooms in on just the spatial cross-validation results.

### Cross-validation results (true $K=3$ )

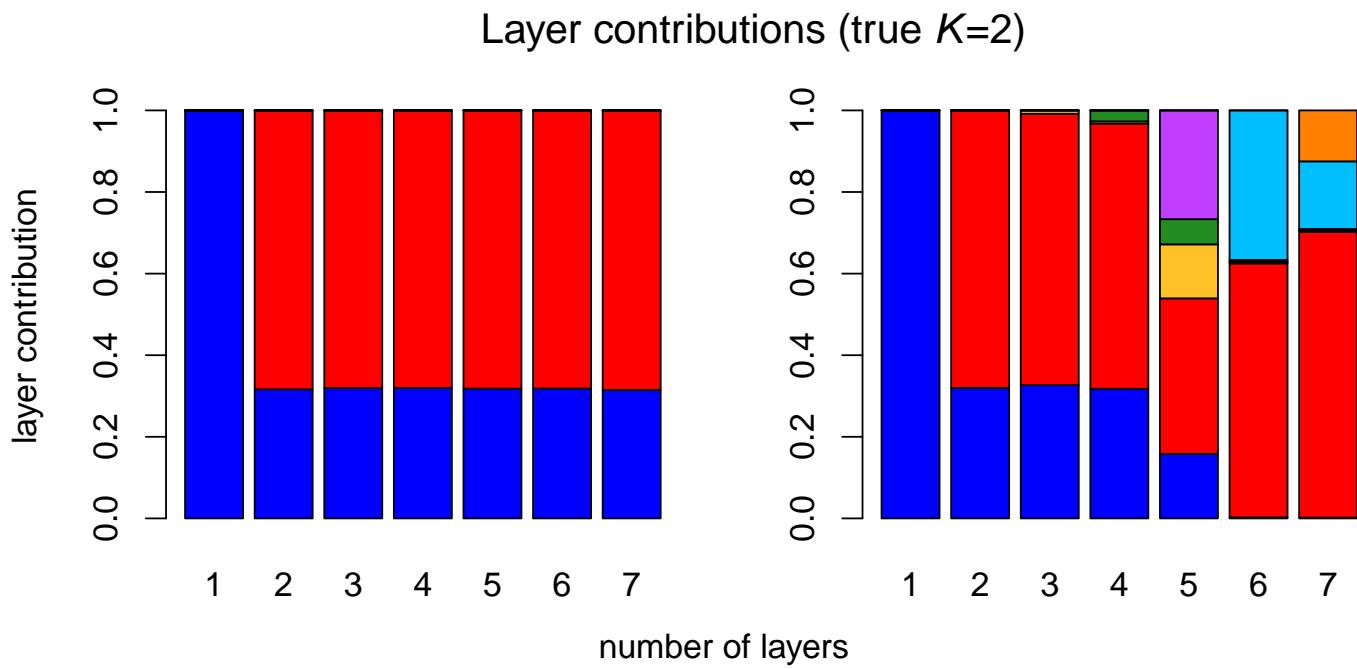


**Figure S14** Cross-validation results for data simulated under  $K = 3$ , comparing the spatial and nonspatial `conStruct` models run with  $K = 1$  through 7. The right panel zooms in on just the spatial cross-validation results.

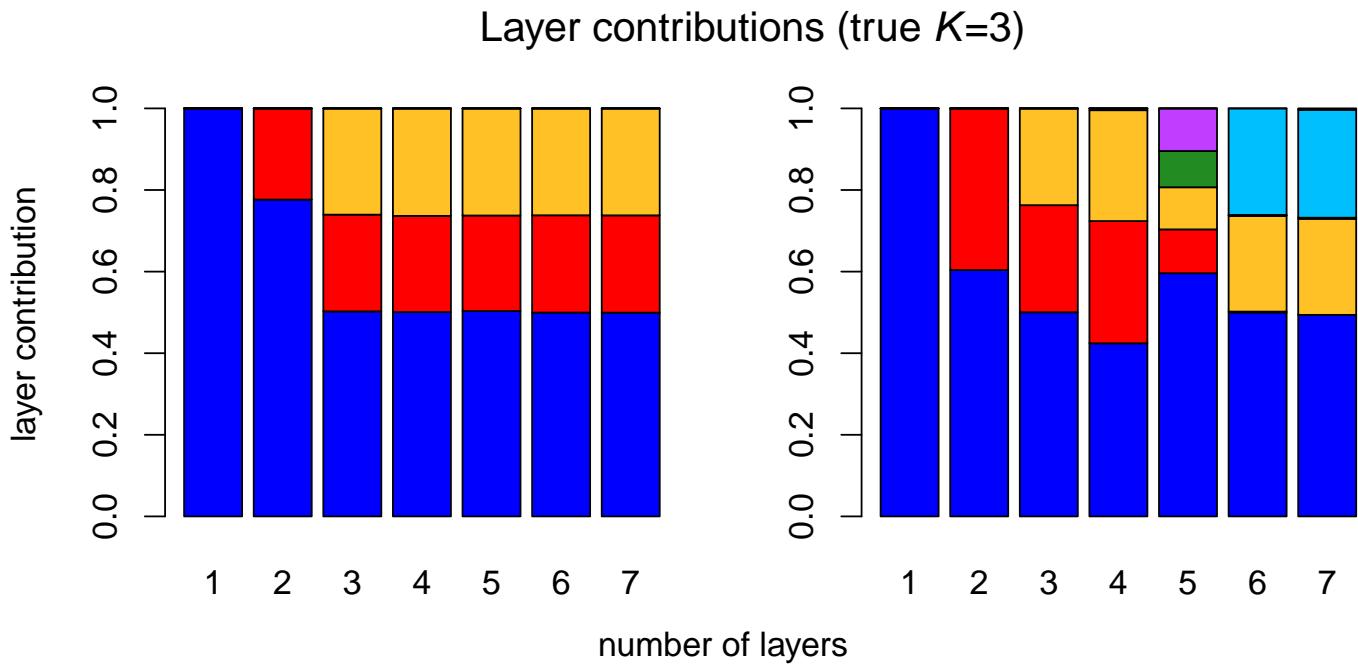
### ADMIXTURE cross-validation results for simulated data



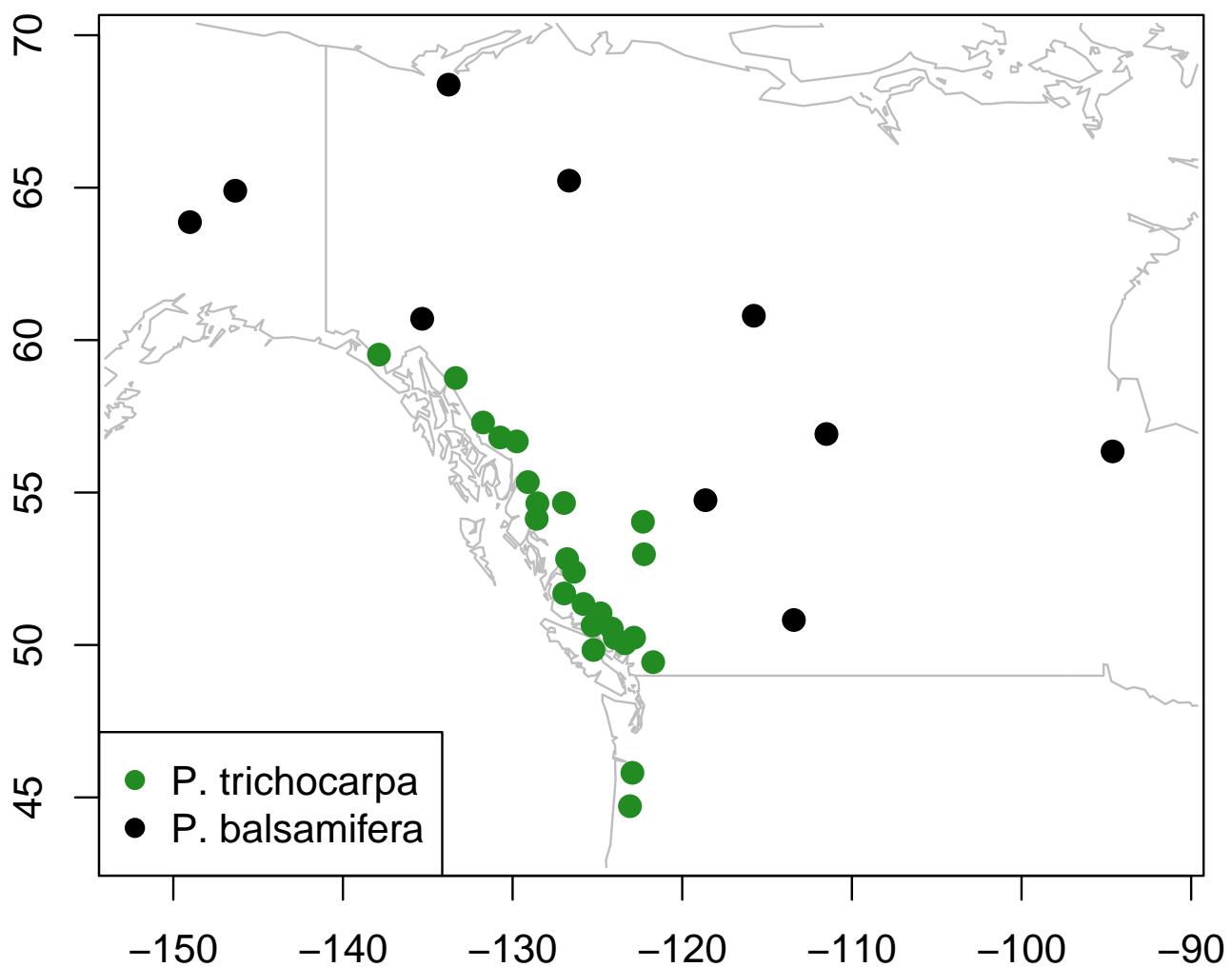
**Figure S15** ADMIXTURE cross-validation results for data simulated under  $K = 1$ ,  $K = 2$ , and  $K = 3$ , run with  $K = 1$  through 7 using 50 data folds (-cv=50). The inset plots zoom in on cross-validation results outlined in the dotted boxes. The preferred model (with the lowest cross-validation error) is highlighted with an asterisk.



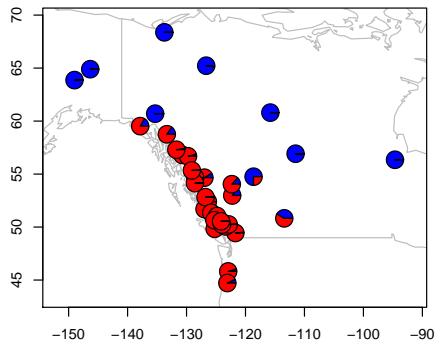
**Figure S16** Layer/cluster contributions (i.e., how much total covariance is contributed by each layer/cluster), for all layers estimated in runs using  $K = 1$  through 7 for the spatial model (left), and for all clusters using the nonspatial `conStruct` model (right). Data were simulated using  $K = 2$ . For each value of  $K$  along the x-axis, there are an equal number of contributions plotted.



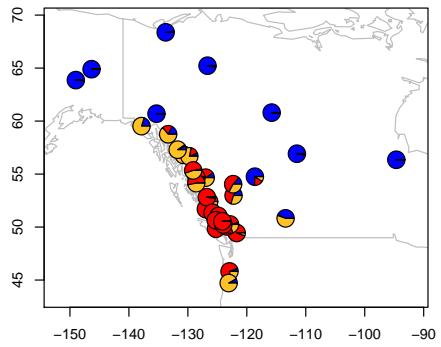
**Figure S17** Layer/cluster contributions (i.e., how much total covariance is contributed by each layer/cluster), for all layers estimated in runs using  $K = 1$  through 7 for the spatial model (left), and for all clusters using the nonspatial `conStruct` model (right). Data were simulated using  $K = 3$ . For each value of  $K$  along the x-axis, there are an equal number of contributions plotted.



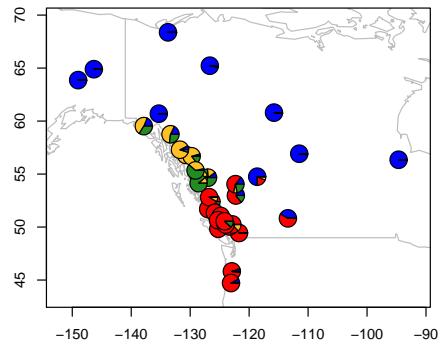
**Figure S18** Map of the sampled *Populus* populations included in the analysis. The color of the sampling location denotes the putative species.



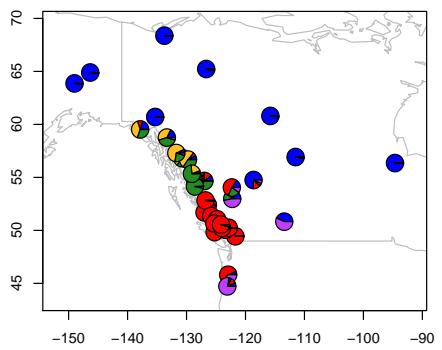
(a)  $K = 2$



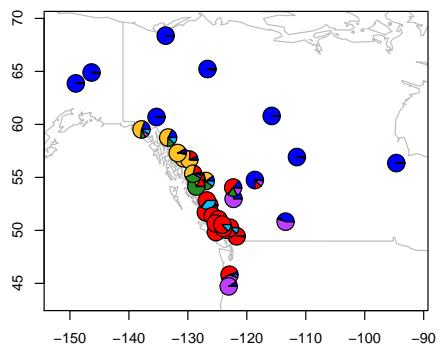
(b)  $K = 3$



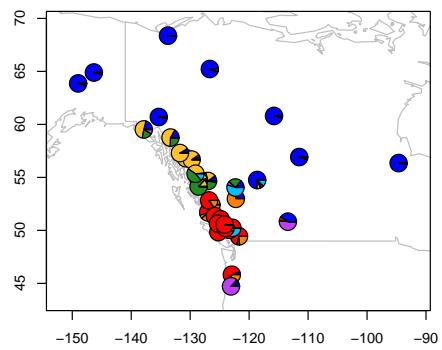
(c)  $K = 4$



(d)  $K = 5$

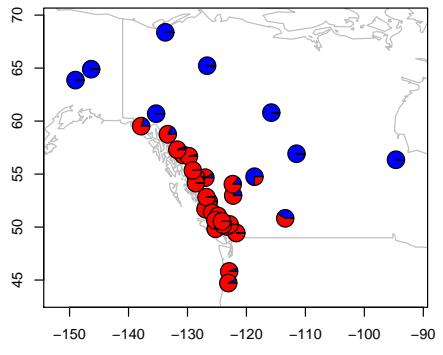


(e)  $K = 6$

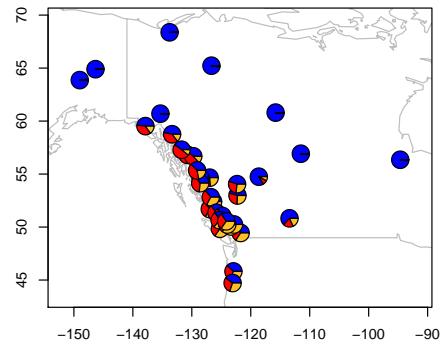


(f)  $K = 7$

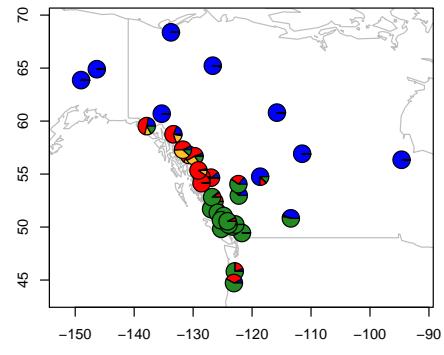
**Figure S19** Maps of admixture proportions estimated for the *Populus* dataset using the spatial conStruct model for  $K = 2$  through 7.



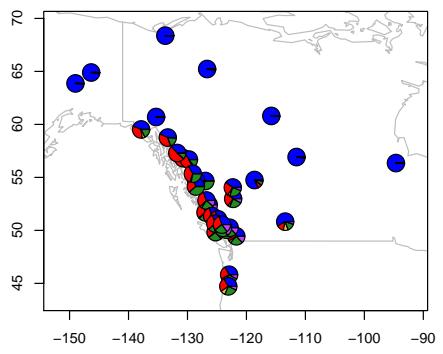
(a)  $K = 2$



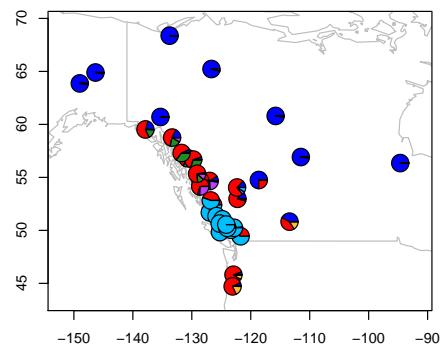
(b)  $K = 3$



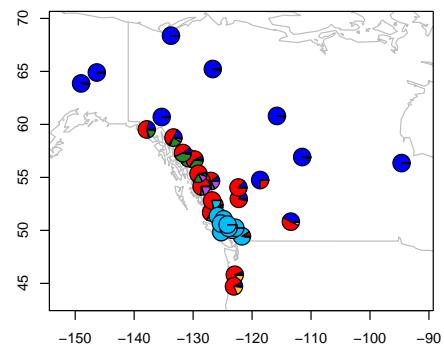
(c)  $K = 4$



(d)  $K = 5$

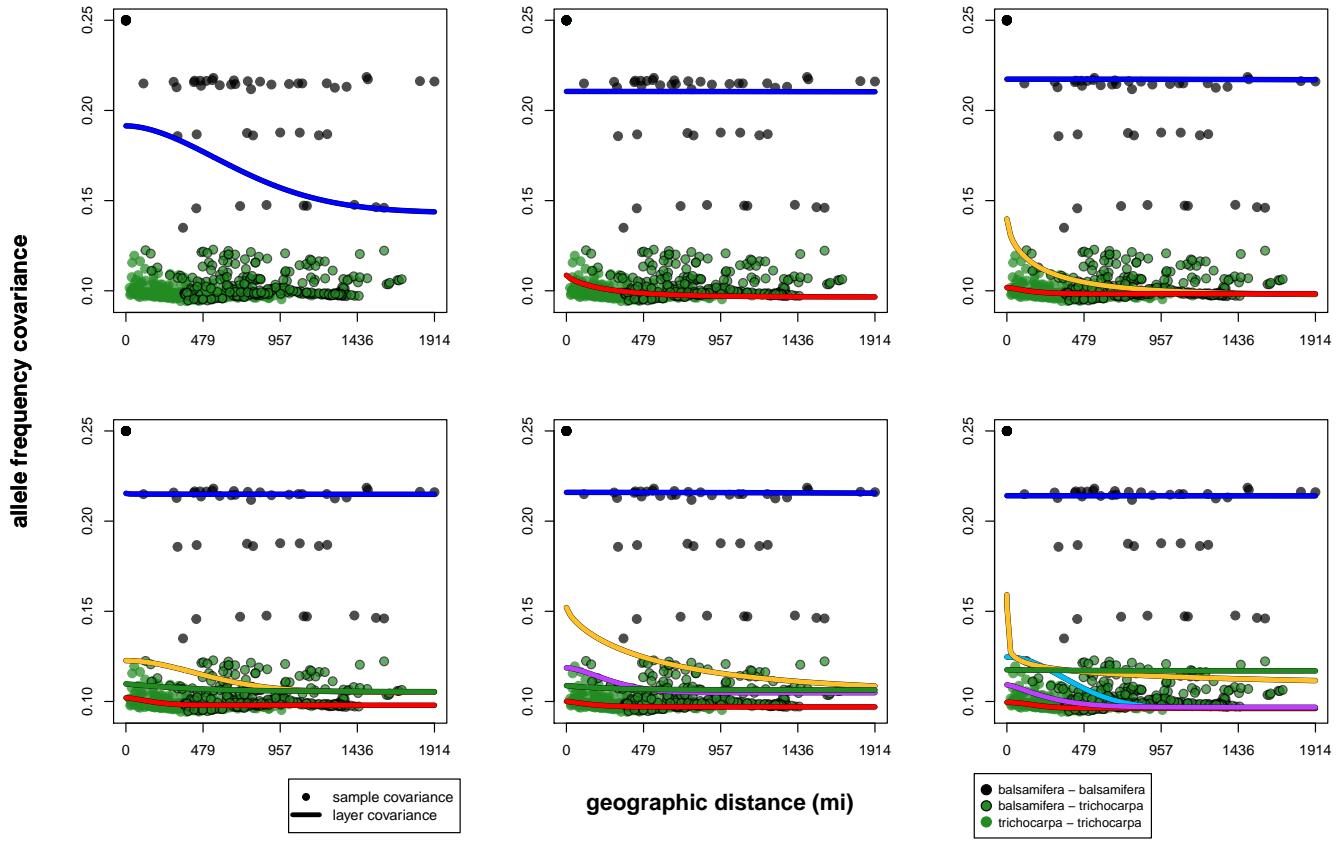


(e)  $K = 6$

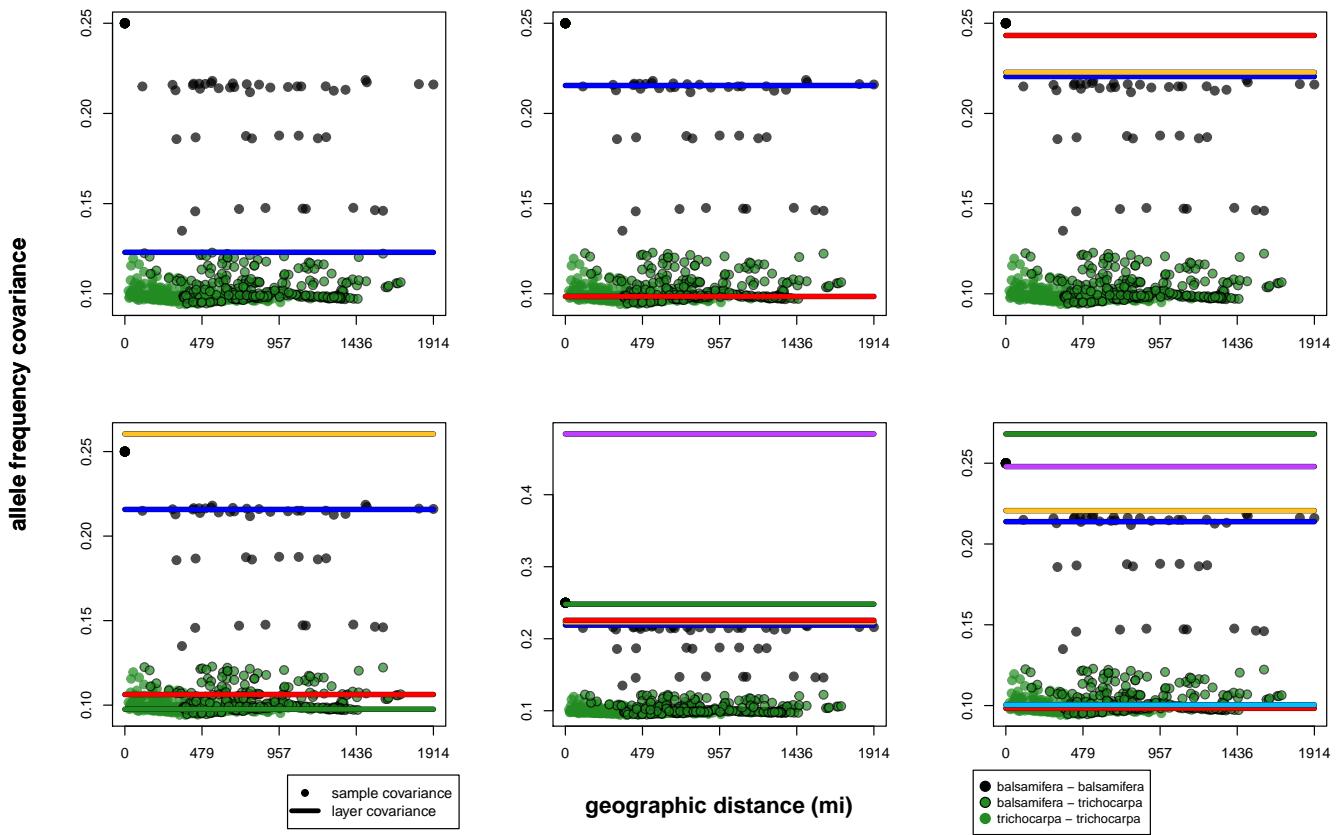


(f)  $K = 7$

**Figure S20** Maps of admixture proportions estimated for the *Populus* dataset using the nonspatial conStruct model for  $K = 2$  through 7.

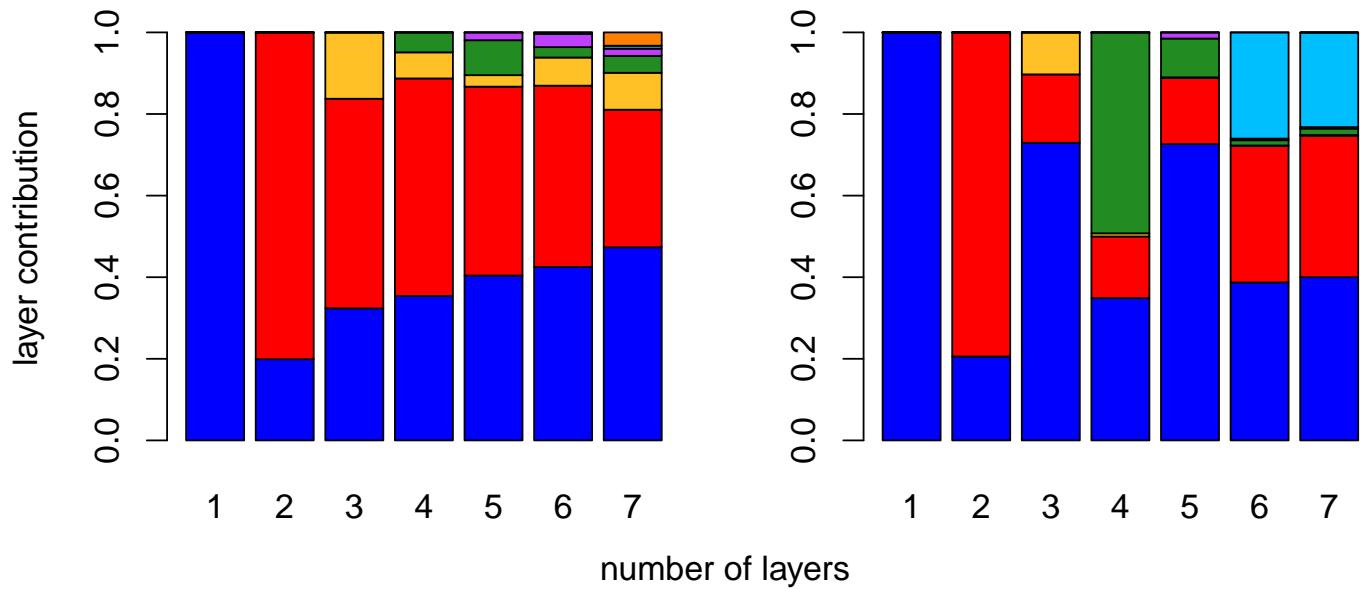


**Figure S21** Plots showing the layer-specific parametric covariance curves estimated for the *Populus* data using the spatial conStruct model run with  $K = 1$  through 6. Line colors are consistent with layer colors in Fig S19. Points are colored by the species they are a covariance between: black on black points are sample covariances between populations of *Populus balsamifera*; green on black points are sample covariances between *balsamifera* and *trichocarpa*; green on green points are sample covariances between *trichocarpa* and *trichocarpa*.

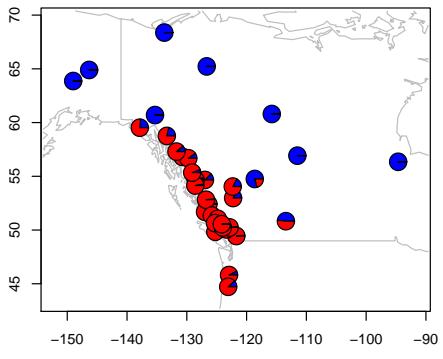


**Figure S22** Plots showing the cluster-specific parametric covariances estimated for the *Populus* data using the nonspatial `conStruct` model run with  $K = 1$  through 6. Line colors are consistent with cluster colors in Fig S20. Points are colored by the species they are a covariance between: black on black points are sample covariances between populations of *Populus balsamifera*; green on black points are sample covariances between *balsamifera* and *trichocarpa*; green on green points are sample covariances between *trichocarpa* and *trichocarpa*.

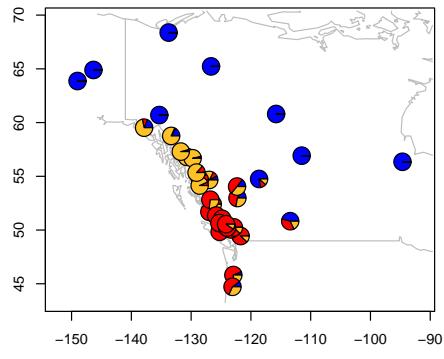
## Layer contributions (Poplars)



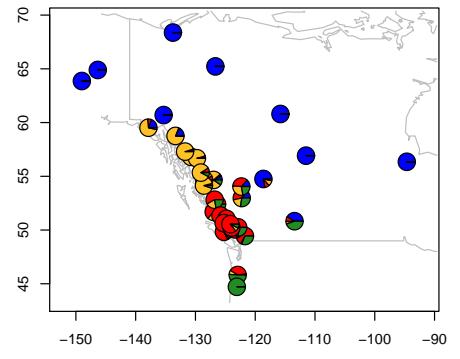
**Figure S23** Layer/cluster contributions (i.e., how much total covariance is contributed by each layer/cluster), for all layers estimated in runs using  $K = 1$  through 7 for the spatial model (left), and for all clusters using the nonspatial `conStruct` model (right). For each value of  $K$  along the x-axis, there are an equal number of contributions plotted. Colors are consistent with Figs S19, S21, S20, and S22.



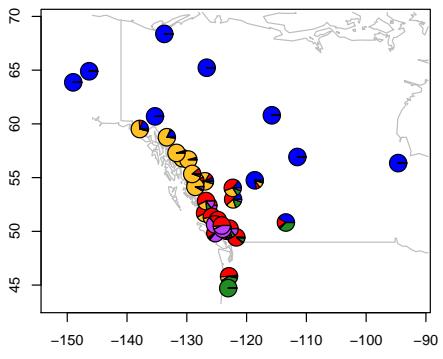
(a)  $K = 2$



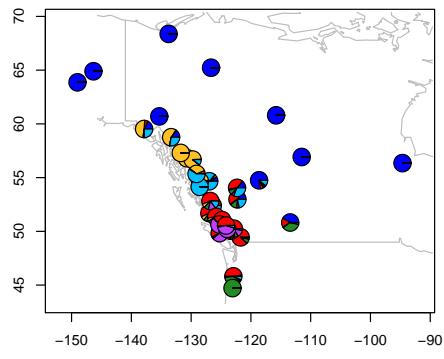
(b)  $K = 3$



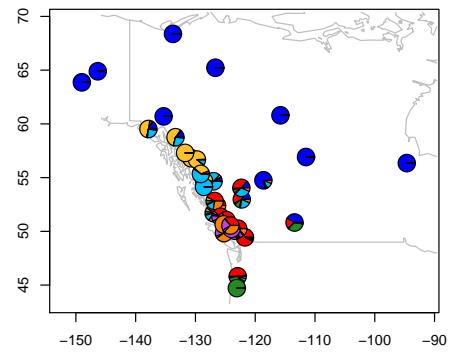
(c)  $K = 4$



(d)  $K = 5$



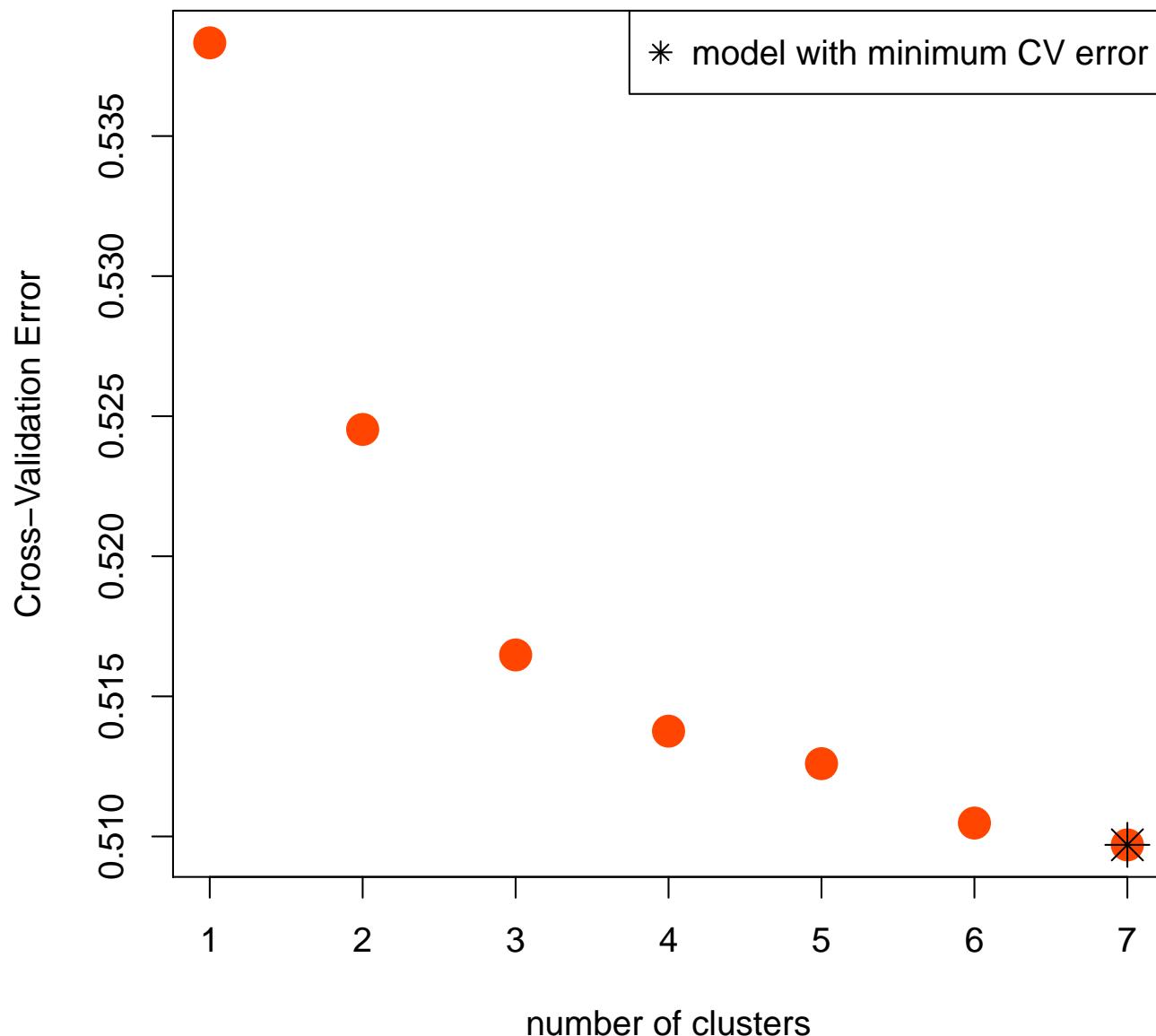
(e)  $K = 6$



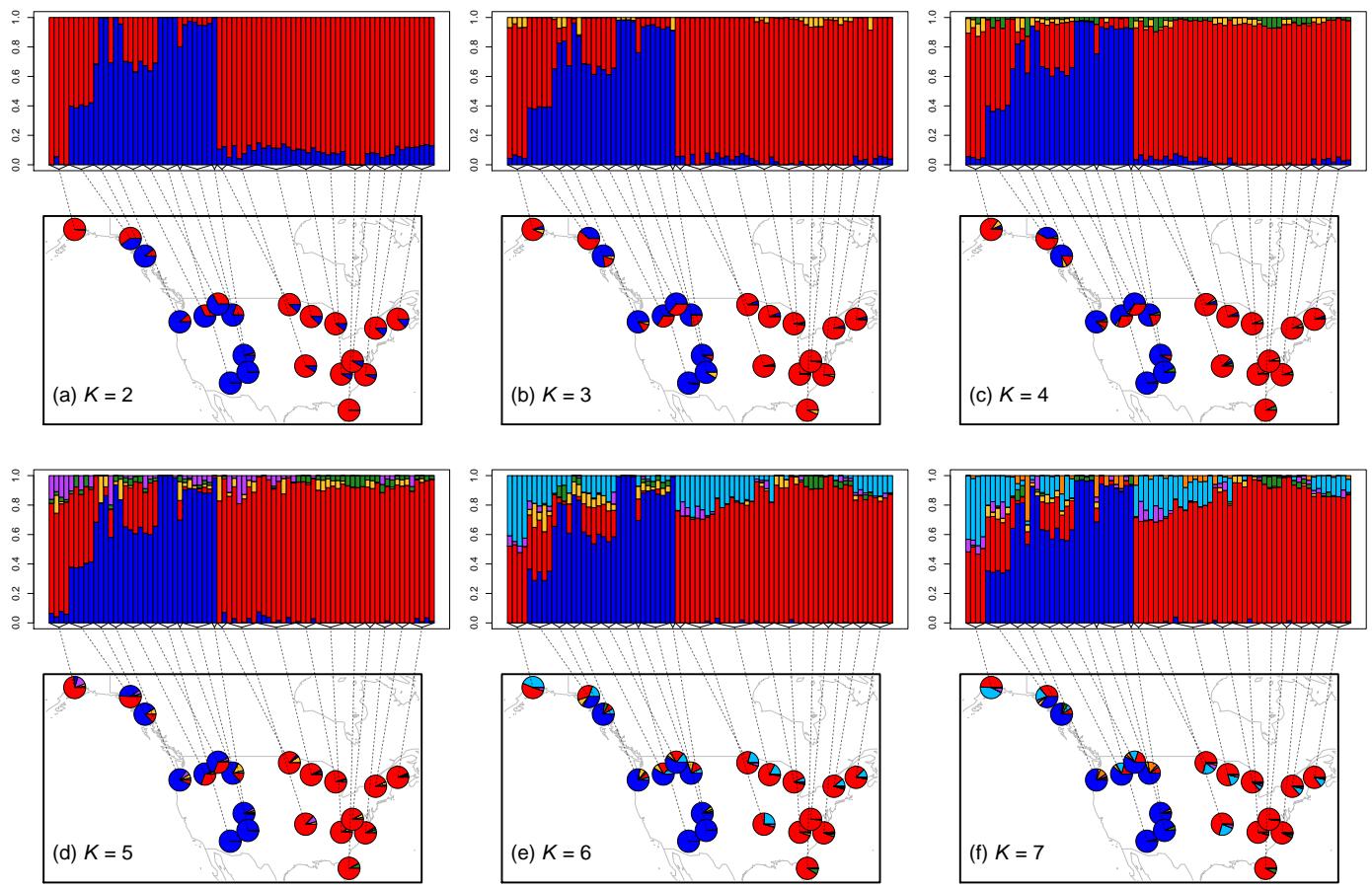
(f)  $K = 7$

**Figure S24** Maps of admixture proportions estimated for the *Populus* dataset using ADMIXTURE [Alexander et al. \(2009\)](#) for  $K = 2$  through 7.

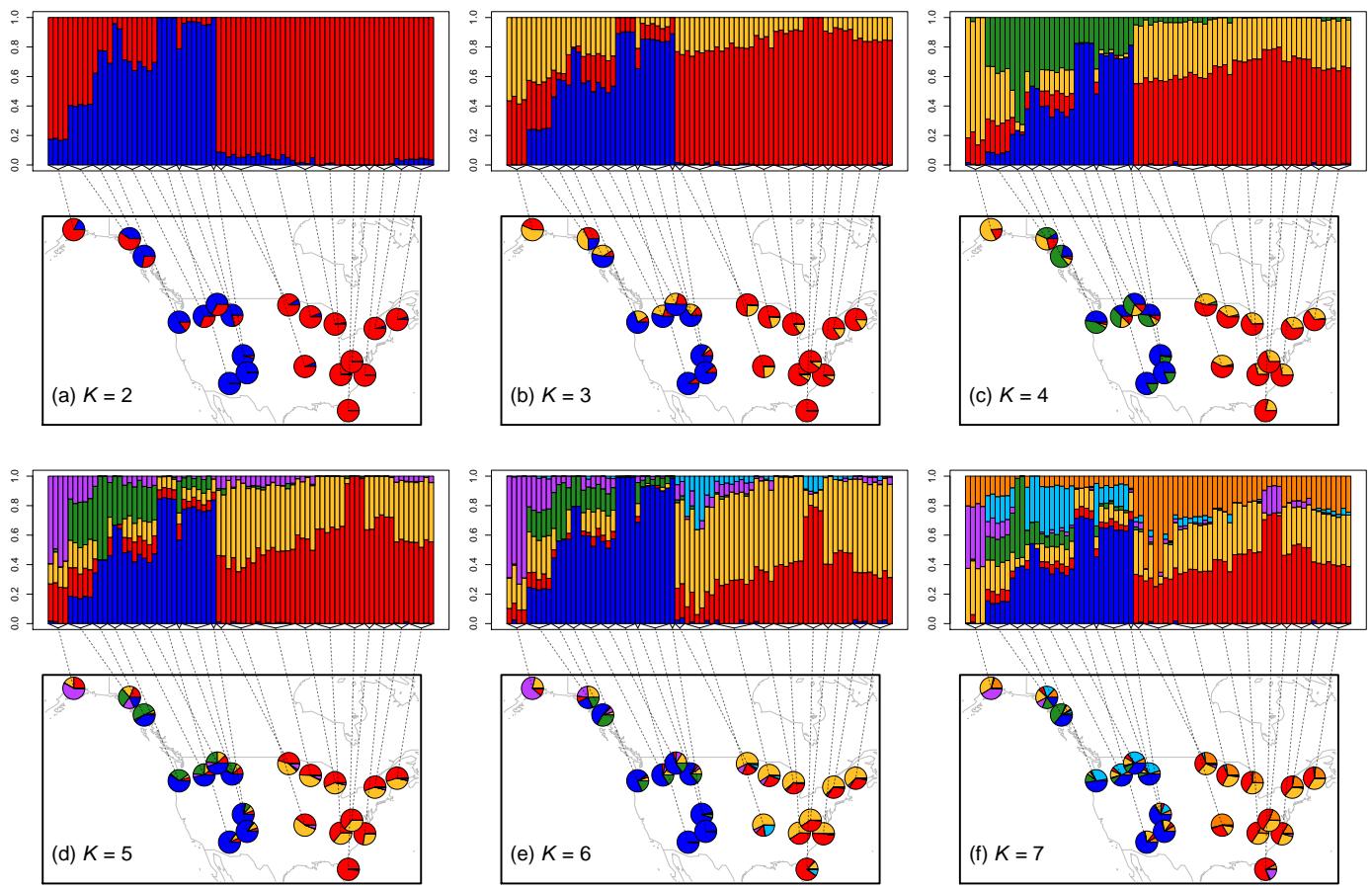
## Poplar ADMIXTURE cross-validation results



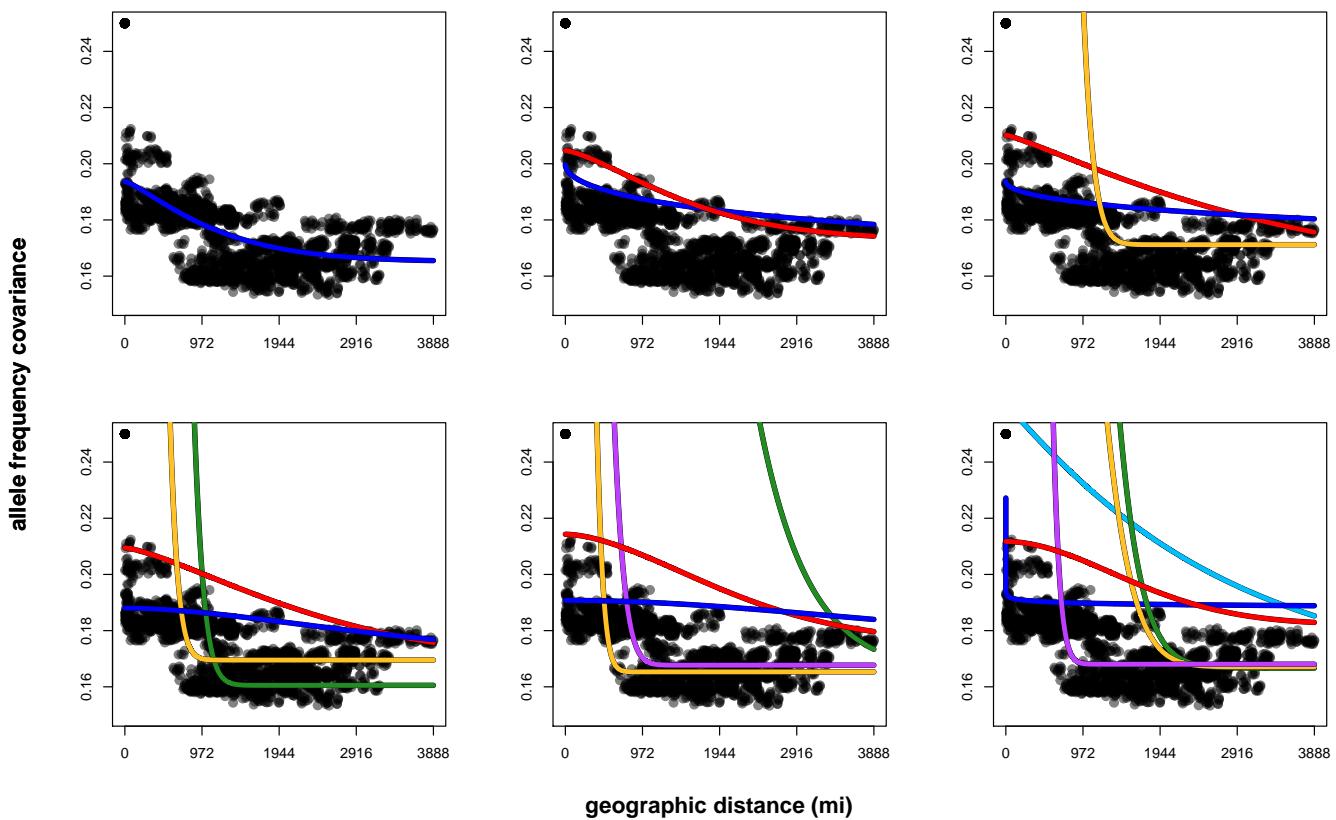
**Figure S25** ADMIXTURE cross-validation results for poplar data, run with  $K = 1$  through 7 using 50 data folds ( $-cv=50$ ). The preferred model (with the lowest cross-validation error) is highlighted with an asterisk.



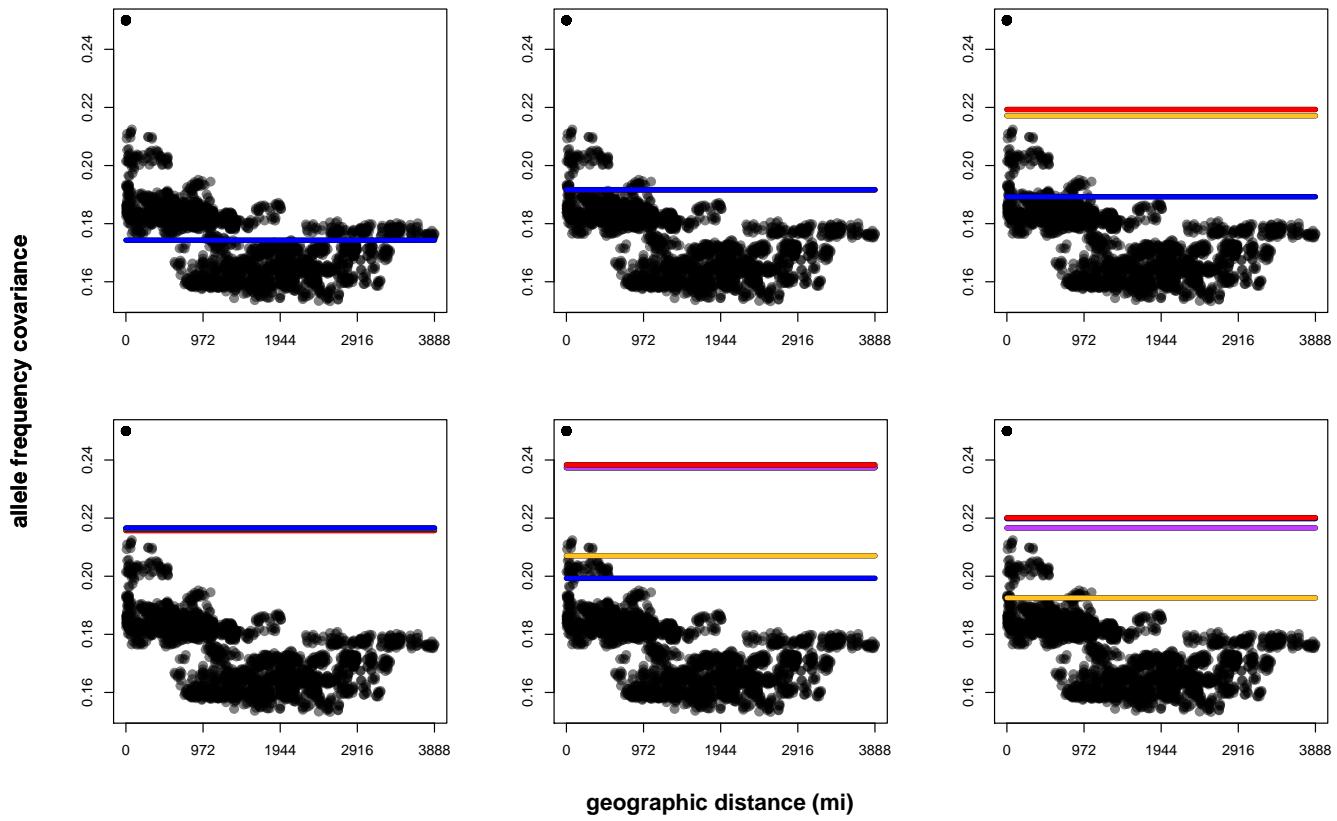
**Figure S26** Map of admixture proportions estimated for the bear dataset using the spatial conStruct model for  $K = 2$  through 7.



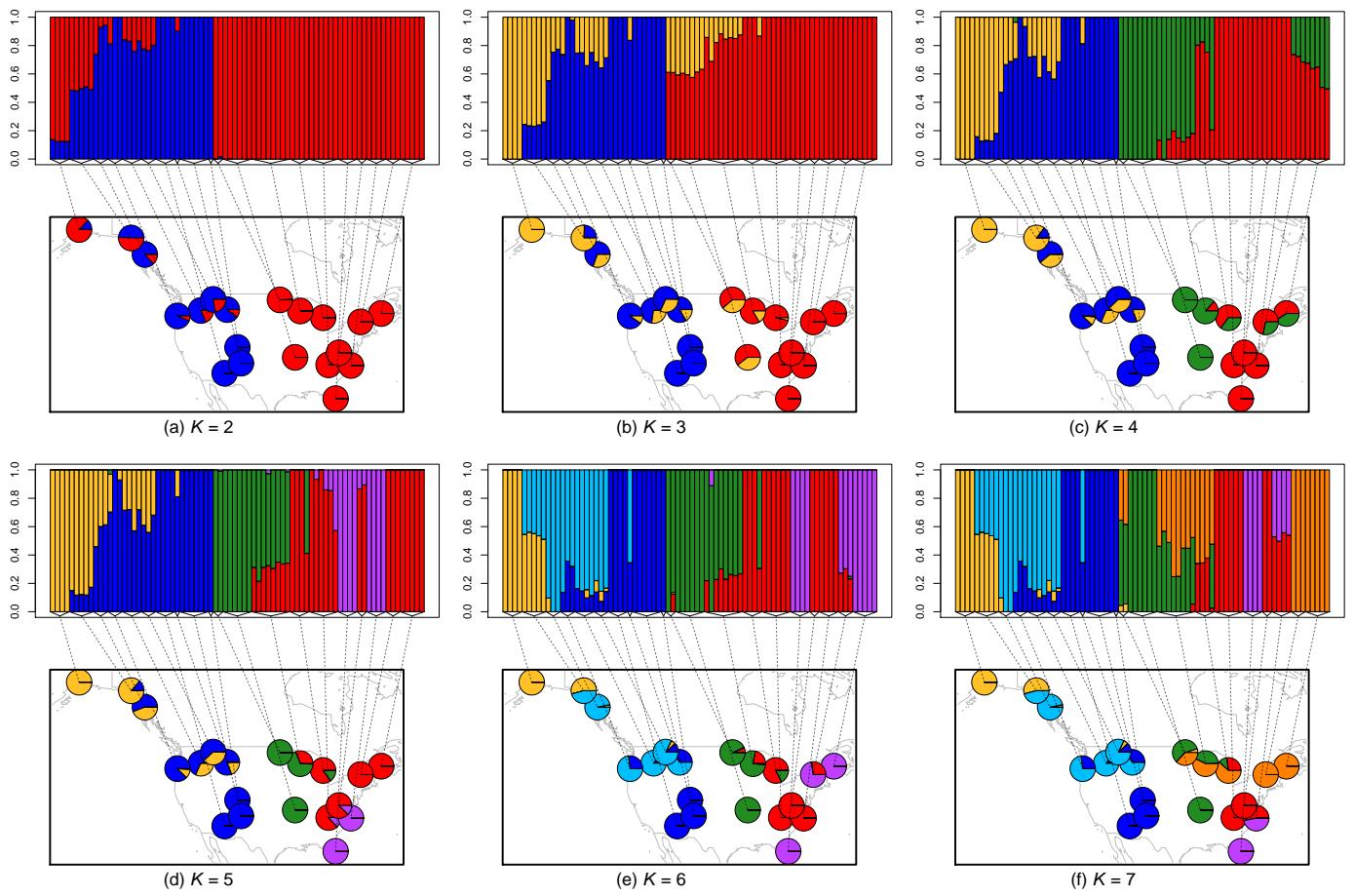
**Figure S27** Map of admixture proportions estimated for the bear dataset using the nonspatial conStruct model for  $K = 2$  through 7.



**Figure S28** Plots showing the layer-specific parametric covariance curves estimated for the black bear data using the spatial conStruct model run with  $K = 1$  through 6.

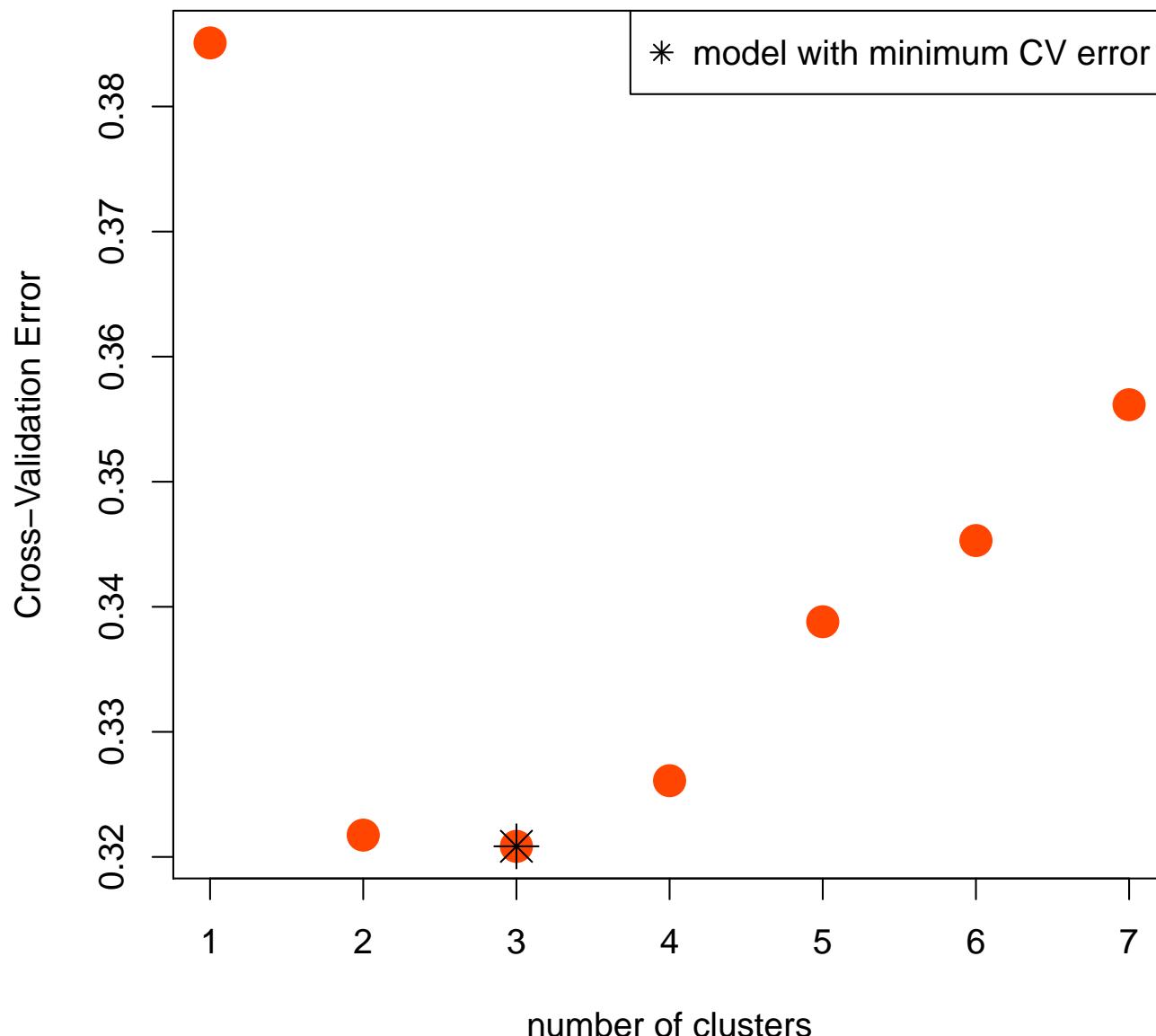


**Figure S29** Plots showing the cluster-specific parametric covariances estimated for the black bear data using the nonspatial conStruct model run with  $K = 1$  through 6.

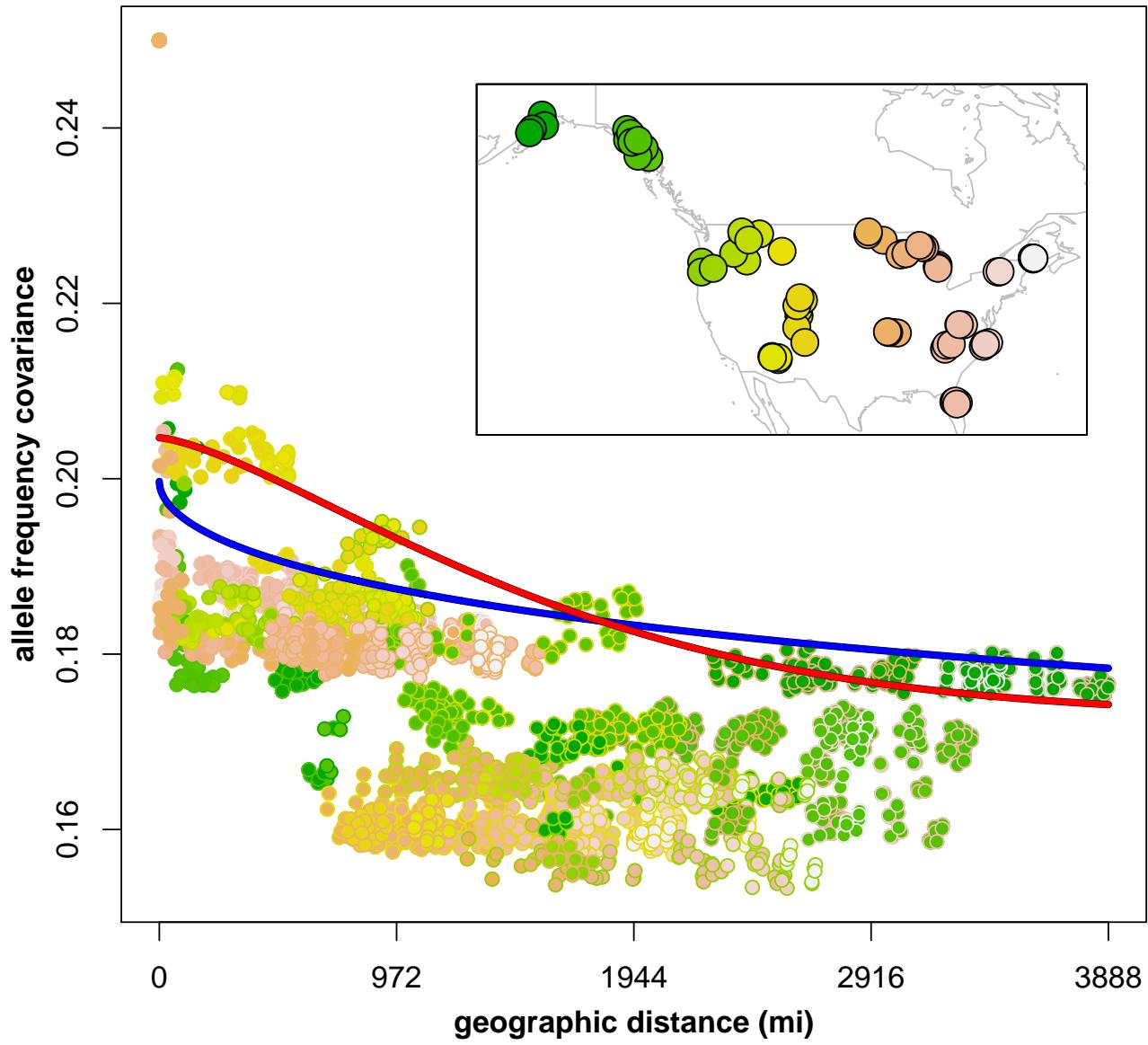


**Figure S30** Maps of admixture proportions estimated for the bear dataset using ADMIXTURE [Alexander et al. \(2009\)](#) for  $K = 2$  through 7.

## Bear ADMIXTURE cross-validation results



**Figure S31** ADMIXTURE cross-validation results for bear data, run with  $K = 1$  through 7 using 50 data folds ( $-cv=50$ ). The preferred model (with the lowest cross-validation error) is highlighted with an asterisk.



**Figure S32** Allele frequency covariance between each pair of sampling locations for bear data, with each point colored to show which two locations it corresponds to (see inset for colors). For instance, set of points showing covariance around 0.18 and above 2000 miles distant have dark green centers with white-to-orange borders, indicating they are comparisons between the Northwest-most (Alaskan) bears and the bears on the Eastern half of the map. The two curves show the decay of covariance within each layer for the spatial conStruct model with  $K = 2$ , as in figure S28.