# A Method for Detecting Population Genetic Structure in Diverse, High Gene-Flow Species

Ryan P. Kelly*, Thomas A. Oliver*, Arjun Sivasundar, and Stephen R. Palumbi

Department of Biology, Hopkins Marine Station of Stanford University, Oceanview Blvd. Pacific Grove, CA 93950.

*These authors contributed equally to this work.

Address correspondence to R. Kelly at the address above, or e-mail: rpk@stanford.edu.

## Abstract

Detecting small amounts of genetic subdivision across geographic space remains a persistent challenge. Often a failure to detect genetic structure is mistaken for evidence of panmixia, when more powerful statistical tests may uncover evidence for subtle geographic differentiation. Such slight subdivision can be demographically and evolutionarily important as well as being critical for management decisions. We introduce here a method, called spatial analysis of shared alleles (SAShA), that detects geographically restricted alleles by comparing the spatial arrangement of allelic co-occurrences with the expectation under panmixia. The approach is allele-based and spatially explicit, eliminating the loss of statistical power that can occur with user-defined populations and statistical averaging within populations. Using simulated data sets generated under a stepping-stone model of gene flow, we show that this method outperforms spatial autocorrelation (SA) and $\Phi_{ST}$ under common real-world conditions: at relatively high migration rates when diversity is moderate or high, especially when sampling is poor. We then use this method to show clear differences in the genetic patterns of 2 nearshore Pacific mollusks, *Tegula funebralis* (= *Chlorostoma funebralis*) and *Katharina tunicata*, whose overall patterns of within-species differentiation are similar according to traditional population genetics analyses. SAShA meaningfully complements $\Phi_{ST}/F_{ST}$, SA, and other existing geographic genetic analyses and is especially appropriate for evaluating species with high gene flow and subtle genetic differentiation.

**Key words:** allele statistic, geographic genetics, population genetics

A central problem in population genetics is detecting small amounts of population subdivision created by substantial, though nonrandom, gene flow among demes. As gene flow increases, values of $F_{ST}$ and its analogues (such as $\Phi_{ST}$) become small relative to their confidence intervals (Waples 1998), making it impossible to assess subtle genetic subdivision without impractically large sample sizes. However, such apparently minor deviations from panmixia can have major demographic and evolutionary implications. For example, identifying distinct stocks of a commercial fish species is critical for management, yet for sufficiently large populations, a migration rate of 10% between distinct stocks may be impossible to distinguish from panmixia (Palumbi 2003). More sensitive statistical tools are clearly desirable.

Previous authors have proposed measures of genetic subdivision based on the observation that alleles arise by mutation and spread over geographic space via migration, making use of allelic identity and the relationships among

alleles to infer limitations to gene flow. Slatkin (1981, 1985) demonstrated that the frequency of alleles occurring in only one geographic location (private alleles) estimates the number of migrants among populations in a given generation (Nm). Subsequently, Hudson (2000) introduced a measure, $S_{nn}$, based not on identical alleles but rather on nearest neighbors in sequence space (i.e., the most similar, nonidentical alleles) that co-occur in a location. He found that the frequency of these co-occurring nearest neighbors could be used to identify significant departures from panmixia, as geographically restricted, newly arisen nearest-neighbor alleles suggest a limitation to the spread of these alleles over space.

Whereas Slatkin's rare alleles method considers only private alleles occurring within a location, and Hudson's nearest-neighbor statistic evaluates closely related alleles occurring within a location, we introduce here a method of analyzing co-occurrences of the same allele over geographic

space. This approach, called spatial analysis of shared alleles (SAShA), is a simple allele statistic that is sensitive to subtle genetic subdivision in ecological time, does not require user-designated populations, is amenable to a variety of data types, and is free from the assumptions required by an underlying theoretical model of gene flow.

Co-occurrences of an allele in different geographic locations are evidence of gene flow, under the assumption that alleles identical in state are identical by descent. We test whether the geographic distances between co-occurrences of alleles are distributed randomly; nonrandom distributions of these pairwise occurrences can indicate departures from panmixia. Alleles may be "underdistributed" (or over-distributed), occurring more closely together (or further apart) in space than expected by chance. In addition to analyzing all alleles in a data set simultaneously, our method allows individual alleles to be scrutinized. This may be useful in cases where common and widespread alleles mask the nonrandom distribution of less common alleles. The statistic returned by SAShA is the average geographic distance between occurrences of alleles, which is intuitive and biologically relevant. Homoplasy, or the existence of alleles that are identical in state but not identical by descent, is expected to occur relatively rarely and without regard to geography and therefore does not introduce systematic bias into the analysis.

Individual- or allele-based methods for detecting geographic genetic structure have become increasingly common in recent years, in part because they address 2 primary drawbacks of traditional population-based statistics: the approximation of biological populations as collection locations and population-level averaging over individuals or alleles. $F_{ST}$ and its relatives (Wright 1951, 1965; Nei 1973; Weir and Cockerham 1984; Excoffier et al. 1992; Slatkin 1995), Fisher's Exact test (Hudson et al. 1992; Raymond and Rousset 1995), the Mantel test (Smouse et al. 1986), Hudson's nearest neighbor (Hudson 2000), and various coalescence-based analyses (see, e.g., Beerli 2006) all suffer a loss of resolution due to these population-level simplifications.

By contrast, SAShA and other allele- or individual-based methods, such as spatial autocorrelation (SA) (Heywood 1991; Hardy and Vekemans 1999; Smouse and Peakall 1999; Epperson 2003), nested clade analysis (Templeton 1998, 2004, 2008; Panchal and Beaumont 2007; Petit 2007), allelic aggregation index analysis (AAIA) (Miller 2005), and "*Sp*" (Vekemans and Hardy 2004) avoid these drawbacks. In addition, such measures make use of geographic information that is ignored in most population statistics (though the Mantel test is an exception), and therefore may detect small degrees of genetic structuring that are missed by population statistics.

Most recently, Novembre and Slatkin (2009) exploited the geographic information content of identical low-frequency alleles to estimate the likelihood distribution for dispersal in 2D space. Like SAShA, this method assumes that alleles arise only once such that identical alleles are treated as identical-by-descent. Unlike our approach,

however, likelihood method of Novembre and Slatkin requires data from individuals randomly sampled across a landscape and a user-specified subset of alleles for analysis. Although these additional requirements make possible a powerful estimation of the size and shape of dispersal, the concept we describe here is far simpler and less computationally intensive, merely revealing the arithmetic mean geographic distance between pairs of shared alleles.

SAShA is expected to be particularly useful for high gene flow species, in which migration occurs much faster than the combined effects of mutation and drift, and shared alleles are consequently spread across the landscape by migration. The geographic distribution of these alleles reflects the spatial extent over which gene flow occurs. Below, we find that the SAShA statistic performs similarly to $\Phi_{ST}$ in modified stepping-stone simulations, but it is more likely than $\Phi_{ST}$ to detect structure when levels of migration are relatively high, when diversity is relatively high, or when sampling is sparse. Both SAShA and $\Phi_{ST}$ substantially outperform spatial autocorrelation in our simulations, though this may be a result of the discrete sampling scheme implemented. We also compare SAShA with another allele-based statistic, Miller's AAIA (see Miller 2005), finding that these 2 perform similarly but that SAShA is not subject to AAIA's high false-negative rate. Finally, we contrast 2 real-world mitochondrial DNA (mtDNA) data sets from the nearshore marine environment to demonstrate the utility of SAShA to uncover genetic structure where other statistics do not.

## Materials and Methods

Our approach requires only a table of allele occurrences by collection localities and the geographic distances between each pair of those localities. Using these data, we compare the observed distribution of geographic distances between instances of each allele (observed distance distribution [ODD]) with a null distribution generated from the same data. This null distribution (expected distance distribution [EDD]) is the distribution of geographic distances between all pairs of samples in the data set regardless of allelic identity and therefore represents the expectation under panmixia (see Appendix 1 for algorithms). We test for significant deviation of the arithmetic mean of the observed distribution from that of the null expectation, reporting the observed mean (OM) distance between co-occurrences of an allele as our statistic.

$$OM = mean(ODD) \qquad (1)$$

When OM is less than the expected mean, the alleles are underdistributed in the aggregate; overdistributed alleles generate an OM greater than the expected mean. OM has the desirable statistical property of being consistent (i.e., approaching a "true" value as more data are considered), though its limits are specific to the data set being analyzed: the maximum mean geographic distance between shared alleles will vary with each data set and field sampling regime. The OM may be conveniently visualized in the context of

the observed and expected distributions as a histogram or cumulative distribution frequency plot, both of which are provided below in the Results section. We note that, although measures of central tendency other than the arithmetic mean might also be useful for distilling the distributions and evaluating the differences between them, the arithmetic mean performs admirably. We explicitly considered a wide variety of such measures in deriving our method, evaluating various nonnormal haplotype frequency distributions, and ultimately selected the arithmetic mean as a simple and powerful way of expressing the statistic.

The significance of the statistic is determined by permutation of the input haplotype matrix, keeping row and column sums constant (see Appendix 1 for algorithm). OM is calculated for each randomly permutated matrix, creating a distribution for the statistic. The *P* value reported for OM is the proportion of those permuted data sets having an OM more extreme than that of the observed data set. For the results presented below, 500 permutations were performed for each data set; larger numbers of permutations result in more precise estimates of significance.

The above calculations may be applied to the entire data set or to any subset thereof. An analysis of each allele individually, for example, may be used to identify individual alleles that drive the pattern observed in the overall data set or else to examine conflict among alleles. Individual allelic distributions are most easily visualized in a cumulative distribution function plot, as below in the Results section. A jackknife procedure, in which the data set is repeatedly reanalyzed after excluding each allele in turn, may also be applied to evaluate the robustness of the pattern observed in the overall data set. We have included basic versions of each of these functions in the source code and compiled executable file, available at http://sasha.stanford.edu.

The performance of the OM statistic was tested in 2 ways: 1) by generating simulated data sets with varying parameter sets and comparing OM with $\Phi_{ST}$, AAIA, and SA and 2) by evaluating 2 real-world mtDNA data sets with SAShA as well as traditional statistics.

### Simulations

The stepping-stone model used to generate simulated data sets was a one locus, many-allele model with nonoverlapping generations implemented with the coalescence-based simulator SIMCOAL (Excoffier et al. 2000). This simulated gene flow along a linear array of 31 equally spaced demes at equilibrium with only the central 9 demes sampled to remove edge effects. Deme sizes were set to remain constant at $10^5$ individuals, sampling either 10 or 25 individuals from each deme for analysis. Mutation rates were set at $2.02 \times 10^{-6}$, $6.05 \times 10^{-6}$, $1.00 \times 10^{-5}$, and $2.02 \times 10^{-5}$ to generate data sets with average $\theta_\pi$ values of 1, 3, 5, and 10. The fragment size generated was 315 bp, mimicking a typical real-world data set of mtDNA haplotypes.

Rather than restricting migration to occur only between immediately adjacent populations, we implemented a Gaussian dispersal kernel, with decreasing probability of migration increasing with distance between populations. We then centered the dispersal kernel on each of the 31 simulated demes, creating the migration matrix input into SIMCOAL, with migration truncated at the edges of the array, representing larvae that dispersed beyond the bounds of the simulation. The width of the dispersal kernel ($\sigma$) was varied to reflect migration out of the focal population at a proportion of $10^{-4}$, $10^{-3}$, $10^{-2}$, $10^{-1}$, and $10^0$. The corresponding $\sigma$ values are 0.225, 0.260, 0.310, and 0.416, with the last case being approximated by a very large value of $\sigma$ (the width of the kernel approaching infinity), resulting in a dispersal matrix of even probabilities (i.e., panmixia).

We performed SAShA analyses on these sampled haplotypes using the MATLAB software coding environment (Mathworks, Inc.); each data set was also exported into Arlequin v3.11(Excoffier et al. 2005) for computation of $\Phi_{ST}$. SA analyses were performed in MATLAB based on the implementation described in Smouse and Peakall (1999) and validated by comparing its outputs with those of GENA-LEX (Peakall and Smouse 2006). We performed SA using 2 different measures of genetic difference among samples: pairwise genetic distance and allelic identity in which the distance matrix consisted of ones (identical allele) and zeros (different allele). Both SA calculations performed similarly in our simulations, and consequently, only the pairwise distance implementation is reported in the Results. Finally, we coded AAIA in MATLAB based on the description in Miller (2005).

The significance of each statistic was determined by permutation: for $\Phi_{ST}$ by permuting individuals (haplotypes) among populations (Excoffier 2000), for SA by permuting the SA correlogram ($T^2$ test) (Smouse and Peakall 1999), and for AAIA and SAShA's OM by permuting as described above (500–1000 permutations were used here for each test).

### Testing Using Original Data Sets

We used 2 cytochrome c oxidase, subunit I (COI) mtDNA data sets from Pacific nearshore species to test the performance of the SAShA method with real-world data. In total, 298 individuals of *Tegula funebralis* (hereafter, *Tegula*), a common intertidal herbivorous snail, were collected from 17 locations between northern Vancouver Island, British Columbia and San Diego, CA, by S.R.P. and R.P.K. between 2003 and 2007. Ninety-four individuals of *Katharina tunicata* (hereafter, *Katharina*), an abundant nearshore chiton, were collected from 9 locations between Kachemak Bay, Alaska and Carmel, California by Douglas J. Eernisse and R.P.K. Genomic DNA was extracted from each species using either commercially available columns (Qiagen DNeasy kits; Qiagen Inc.) or chelex (10% in water, incubated at 65 degrees for 1 h). An approximately 650 bp fragment of the COI gene was amplified via polymerase chain reaction with the primers LCO1490 and HCO2198 (Folmer et al. 1994), sequenced using Big Dye version 1.1 (Applied Biosystems, Inc.), and read on an ABI 3700 or ABI 3730xl sequence analyzer. Resulting sequences were then trimmed to 492 bp (*Tegula*) and 357 bp (*Katharina*) to eliminate all
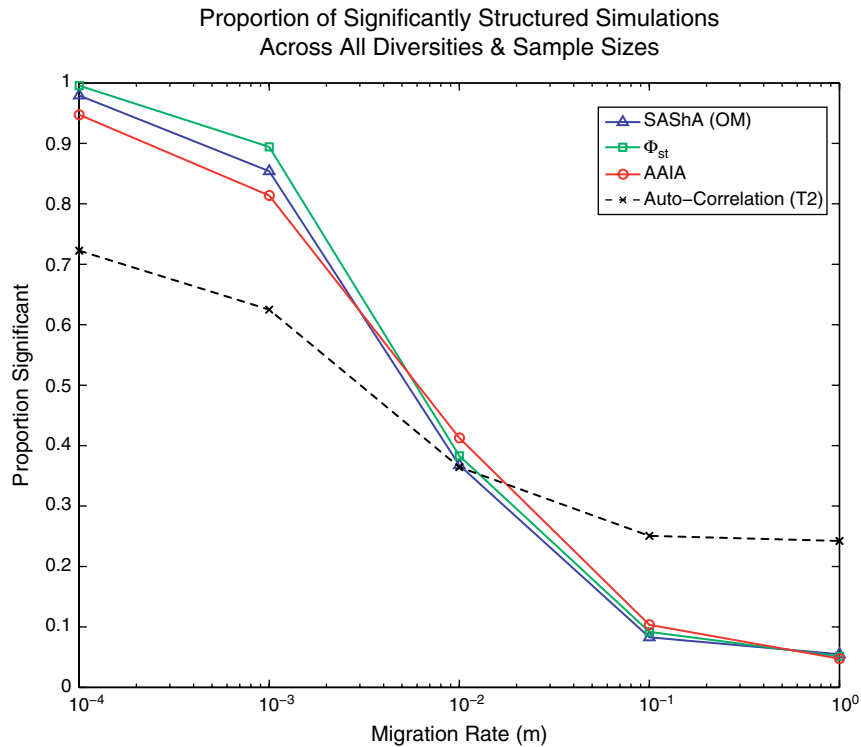
**Figure 1.** The proportion of data sets in which significant genetic structure was detected by SA, Φ, AAIA, and SAShA's OM for 40 000 simulated data sets with varying levels of diversity, shown over the proportion of migrants among demes in a linear stepping-stone model with a Gaussian dispersal kernel.

missing data and input to Arlequin 3.11 in order to calculate $\Phi_{ST}$ and Mantel values and to output the table of haplotype frequencies by population. This table and a matrix of pairwise geographic distances between collection locations were then input into SAShA for analysis.

## Results

### Performance of SAShA with Simulated Data Sets

SAShA and other analyses were carried out on ~40 000 data sets generated using a stepping-stone migration model in SIMCOAL, as detailed in the Materials and Methods. One thousand data sets were generated for each of 40 unique parameter combinations corresponding to each of 5 migration rates ($m = 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}$, and nearly 1), 4 diversity levels ($\theta_\pi = 1, 3, 5$, and 10), and 2 levels of sampling intensity ($n = 10$ per deme or 25 per deme). A small number of SIMCOAL simulations did not converge, and therefore the actual total number of simulations was 39 940.

We compared the performance of SAShA's OM statistic with 3 other population genetics methods, $\Phi_{ST}$, SA, and Miller's AAIA. $\Phi_{ST}$, which calculates the amount of genetic variance in a data set attributable to between-population differences (Excoffier et al. 1992), is a good benchmark for the performance of a statistic because it is

commonly used and its behavior is well understood. SA, although less widely applied by population geneticists, is similar in concept to SAShA and is more often used for landscape genetics with continuous sampling schemes (Manel et al. 2003). It is included here as an allele statistic comparable with SAShA. We finally included a comparison of SAShA with AAIA because the 2 are philosophically similar, both measuring the geographic extent of identically alleles. However, the 2 methods differ in their accounting of that spatial extent: whereas AAIA uses only the nearest-neighbor distance, SAShA encompasses the distribution of all pairs of identical alleles.

We compared the performance of the statistics by reporting the proportion of data sets in which significant genetic subdivision is detected. At the lowest migration level ($m = 10^{-4}$), essentially all data sets were subdivided; those in which no structure is detected are understood to be false negatives (type II errors). Conversely, data sets at the highest migration rate ($m \approx 1$) were panmictic, having equally probable migration between any 2 populations, and thus reflect the false-positive rate (type I error) for each statistic.

Over all parameter sets, the results for all statistics had qualitatively similar sigmoidal shapes, detecting significant structure in the vast majority of data sets with low migration rates ($m = 10^{-4}, 10^{-3}$) and identifying structure less often as migration increases toward panmixia
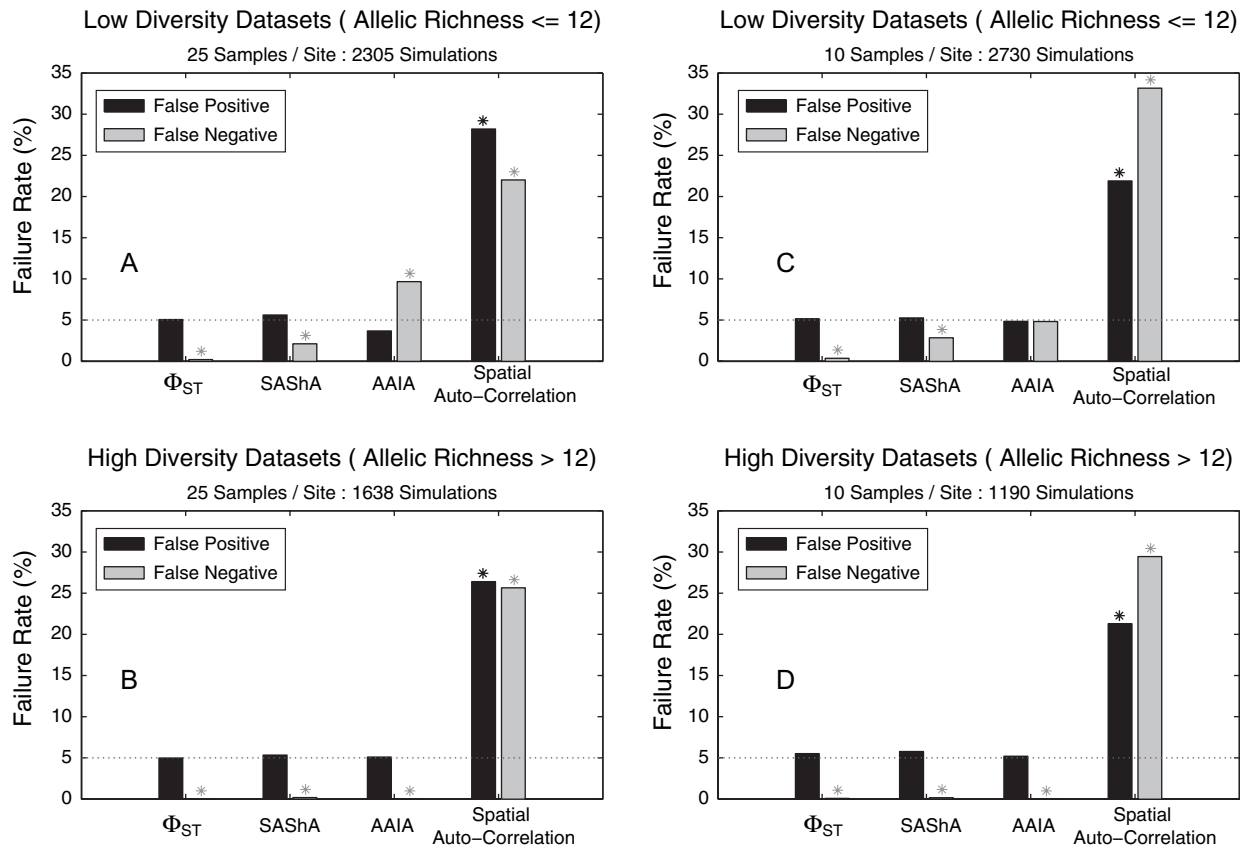
**Figure 2.**   False-positive and false-negative rates for each of the 4 statistics employed: SA, $\Phi_{ST}$, AAIA, and SAShA's OM. Data are grouped by sampling intensity (Figure 2A,B, 25 individuals sampled per locality; Figure 2C,D, 10 individuals sampled per locality), and simulated diversity (Figure 2A,C: $\theta_\pi = 1,3$; Figure 2B,D: $\theta_\pi = 5,10$). Asterisks indicate a significant deviation from a rate of 5% ($\chi^2$ test, bonferroni corrected $\alpha = 0.05/32$ tests).

(Figure 1). SA (9 equally spaced distance bins, corresponding to the 9 sampled simulated populations) detected structure substantially less often than the other 2 methods at low migration, yet maintained strikingly high false-positive and false-negative rates (Figures 1 and 2). This behavior may be due to the discrete sampling scheme employed here: we aimed to simulate how population genetics is often performed in the field, sampling multiple individuals at discrete spatial intervals rather than sampling one or few individuals over a continuous landscape. A thorough examination of the behavior of SA under various sampling regimes is outside the scope of this paper.

SAShA's OM and AAIA show similar results across a broad swath of parameter space; however, AAIA is subject to a high false-negative rate (Figures 1 and 2A,B). This is most likely due to AAIA's use of nearest-neighbor distance to measure the spatial extent of alleles: when sampling is high and migration is low, identical alleles tend to co-occur in the same sampling location, creating a common nearest-neighbor distance of zero. When zero-distance occurrences are very common, they likely swamp out the geographic signal from the other shared alleles across the simulated geography.

This is particularly problematic when allelic richness is low, which makes repeated sampling of the same allele more likely. In data sets with low diversities ($\theta_\pi = 1, 3$) and high sample sizes (25 samples per locality), AAIA fails to detect significant structure in 11.9% of simulations compared with SAShA and $\Phi_{ST}$'s rates of 2.7% and 0.04% for the same data sets (Figure 2A). As SAShA is calculated using the entire distribution of distances between identical alleles, it avoids this pitfall, returning false-negative rates at or below 5% under all conditions (Figure 2).

Regardless of sampling intensity, SAShA consistently performed better in data sets with greater amounts of genetic diversity (Figure 3). At intermediate to high levels of gene flow ($m = 10^{-3}, 10^{-2}$), SAShA detected structure significantly more often in higher diversity data sets ($\theta_\pi = 5, 10$) than in those with lower diversities ($\theta_\pi = 1, 3$; 2-tailed $\chi^2$ test, $P < 0.0001$). At high levels of gene flow ($m = 10^{-2}$) SAShA's OM either shows no significant difference from or significantly outperforms $\Phi_{ST}$ in both higher diversity data sets ($\theta_\pi = 5, 10$) (Figure 3C,D,G,H). The converse is true in lower diversity data sets (Figure 3A,D,E,F).

Reducing the sample size from 25 per deme to 10 erodes the efficacy of both statistics (Figure 3A–D vs. Figure 3E–H).
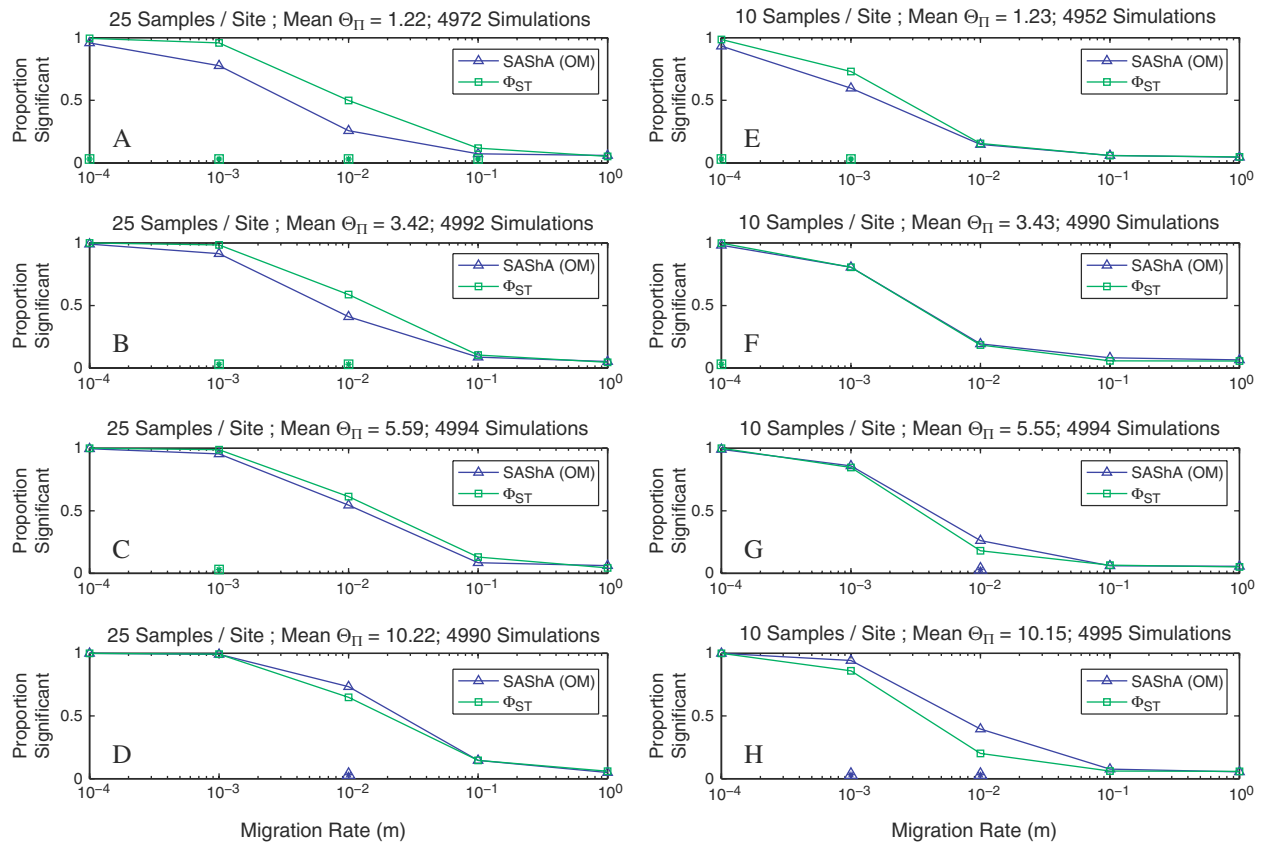
**Figure 3.** The effect of diversity on the power of $\Phi_{ST}$ and SAShA's OM to detect structure in simulations at various migration rates and at different sampling intensities (**A–D**, 25 individuals sampled per locality; **E–H**, 10 individuals sampled per locality). Simulations for each level of simulated diversity are shown (A, E: $\theta_\pi = 1$; **B, F**: $\theta_\pi = 3$; **C, G**: $\theta_\pi = 5$; **D, H**: $\theta_\pi = 10$). A symbol on the $x$ axis indicates a statistically significant difference in power between the 2 statistics shown (Fisher's Exact test, bonferroni corrected $\alpha = 0.05/40$ tests).

However, SAShA's OM statistic is more robust to sparse sampling than $\Phi_{ST}$, showing a less severe decline in its ability to detect structure. OM's relative advantage consequently increases in sparsely sampled data sets with high allelic diversity, and its relative disadvantage at low diversity is substantially reduced (Figure 3E–H).

### Performance of SAShA with Original Data sets

The hidden difference between the arrangement of shared haplotypes in *T. funebralis* and *K. tunicata* further demonstrates the utility of SAShA using biological data. These data sets are typical of the way in which sampling is often done in the field: the samples are of uneven sizes and taken at uneven spatial intervals (Table 1). The 2 species have similar development and larval durations (Strathmann 1987; Moran 1997) and are found in sympatry over much of their ranges. *Tegula* exhibits low genetic variance among populations ($\Phi_{ST} = 0.006$; nonsignificant), an average pairwise difference ($\pi$) of 2.47 among COI haplotypes, and $H = 0.8893$. *Katharina* is comparable by these measurements: $\Phi_{ST} = 0.0218$ (nonsignificant), $\pi = 2.89$, $H = 0.9354$ (Table 1).

The SAShA analysis substantiates *Tegula*'s lack of genetic structure: the spatial arrangement of COI haplotypes is not statistically different from the expectation under panmixia (OM = 778.19 km, expected 756.9 km, nonsignificant; Figure 4A and Table 1). By contrast, *Katharina*'s haplotypes are significantly underdistributed according to the SAShA statistics (OM = 615 km, expected 920 km; *P* value = 0.001; Figure 3B and Table 1), revealing a level of population genetic structure that was otherwise undetected.

Note that although both *Tegula* and *Katharina* have significant Mantel correlations, only the latter has a significant OM. On closer inspection, *Katharina*'s Mantel plot shows a continuous and gradual increase of genetic distance with geographic distance. *Tegula*'s Mantel result is driven by a single data point: the most geographically distant population pair has a positive $F_{ST}$, and removing this comparison eliminates the correlation and its significance (data not shown). This provides a fortuitous test for SASHA: whereas the exceptional population pair drives *Tegula*'s Mantel result, it appears that SAShA is robust to this outlier.

**Table I** Traditional population genetics statistics for *Katharina tunicata* and *Tegula funebralis* based on mtDNA COI data

|  | *K. tunicata* | *T. funebralis* |
| --- | --- | --- |
| No. of locations | 9 | 17 |
| Average no./location | 10.56 | 17.53 |
| No. of individuals sampled | 104 | 298 |
| No. of alleles | 39 | 84 |
| Haplotypic diversity ($H$) | 0.9354 | 0.8893 |
| Maximum geographic Distance sampled (km) | 3324 | 2312 |
| Traditional population genetics results |  |  |
| $\Phi_{ST}$ ($P$ value) | 0.0218 (0.12) | 0.00599 (0.19) |
| $\pi$ | 2.47 | 2.89 |
| $\theta_S$ | 3.02 | 2.56 |
| Tajima's D | $-1.25$ | $-0.08$ |
| Mantel correlation ($P$ value) | 0.719 (0) | 0.279 (0.01) |
| Exact test $P$ value | 1 | 1 |
| SAShA results |  |  |
| Overall OM | 615.49 km | 778.19 km |
| Overall expected mean | 920.63 km | 756.9 km |
| OM $P$ value | 0.001 | 0.4 |

The 2 species appear similar in every respect according to traditional population genetics statistics, with nonsignificant $\Phi_{ST}$ values and exact tests, and significant Mantel correlations. However, the SAShA statistics reveal clear differences between the species: whereas *Tegula*'s haplotypes are not distributed significantly differently from the expectation under panmixia, *Katharina*'s haplotypes are significantly underdistributed.

The haplotype-by-haplotype analysis (Figure A1) makes the differences between the 2 species clearer. Whereas *Tegula*'s haplotypes are all nearly uniformly distributed, *Katharina*'s haplotypes are nearly all underdistributed (number 3 significantly so; $P = 0.001$), driving the pattern seen in the overall OM value. Both species show minor discord among rare haplotypes; in general this is expected—alleles with few instances will not be uniformly distributed by definition—and rare alleles tend to be nonsignificant for this reason. Jackknife SAShA analyses (Table 2) indicate the robustness of the overall results, which remain qualitatively the same when any one haplotype is removed.

## Discussion

We find that the SAShA method consistently reveals genetic structure that is missed by existing genetic analyses, although avoiding the high false-negative rates exhibited by a conceptually similar allele-based statistic. In particular, the approach is effective when diversity is moderate to high ($\theta_\pi = 5, 10$), as is characteristic of many population genetics data sets. The method is effective at relatively high levels of migration, consistent with our expectation that co-occurrences of alleles in space are informative when migration occurs much faster than mutation and drift. Moreover, SAShA returns the calculated geographic distance by which alleles are over- or underdispersed relative to the expectation under panmixia, information that is immediately useful in determining the

spatial scales over which gene flow is observed. Finally, we note that SAShA's relative advantage over $\Phi_{ST}$ is maximized at small sample sizes; such suboptimal sampling is often all that is possible in studies of rare or endangered species, ancient DNA, or difficult-to-obtain individuals from the field.

By using allelic identities rather calculated values of relatedness among pairs of individuals (as reviewed in Vekemans and Hardy 2004), SAShA is subject to the valid criticism that it uses only a subset of the available information in a genetic data set. However, the results above suggest that focusing on the most informative fraction of the available data (i.e., identical alleles) can yield a simple but powerful statistic. This increase in power may be due to minimizing the "noise" introduced into other genetic analyses by estimates of relatedness, by large and variable genetic distances among individuals, and by reliance on an underlying model of gene flow. Alleles identical in state but not by descent (i.e., homoplasy), as discussed further below, diminish SAShA's power but do not introduce systematic bias to the analysis.

### Simulated Data Sets

The results of the SAShA analysis on the simulated data sets demonstrate the advantageous behavior of allele-based statistics in the "Waples zone" of high migration and concomitantly weak structure (see Waples 1998). When migration is low, (Phi)ST, AAIA, and SAShA indicate significant structure in the vast majority of data sets. As migration increases, a shrinking proportion of data sets is significant. This proportion decreases at a different rate for each statistic, depending on the allelic diversity and the sampling intensity. At $m = 10^{-3}$ and $10^{-2}$, SAShA's OM and AAIA detect significantly more structure than $\Phi_{ST}$ when diversity is high ($\theta_\pi = 5, 10$) and sampling is sparse ($n = 10$ per deme). We note that this is an intermediate value for migration in our simulated data sets, but it represents high migration relative to what can be detected confidently with genetic methods in wild populations (Waples 1998; Palumbi 2003).

Both OM and AAIA perform well when diversity is high. However, because AAIA relies on nearest-neighbor distance between identical alleles, when diversity is low and/or when sampling intensity in each locality is high, AAIA dramatically loses power and fails to detect significant structure over 10% of the time (Figure 2A). Note that as AAIA's false positive rate therefore increases as more data are gathered. By employing the distributions of distances between all identical alleles, not a single nearest-neighbor distance, SAShA's OM avoids this pitfall, returning low false-negative rates (Figure 2).

We note that, in simulating deme sizes of $10^5$, we have assumed robust effective population sizes, perhaps larger than might occur in some natural species (especially imperiled species). In general, smaller deme sizes lead to lower diversity and more drift and require higher migration rates (though a smaller absolute number of migrants) to homogenize the demes. Because SAShA's comparative advantage is at high
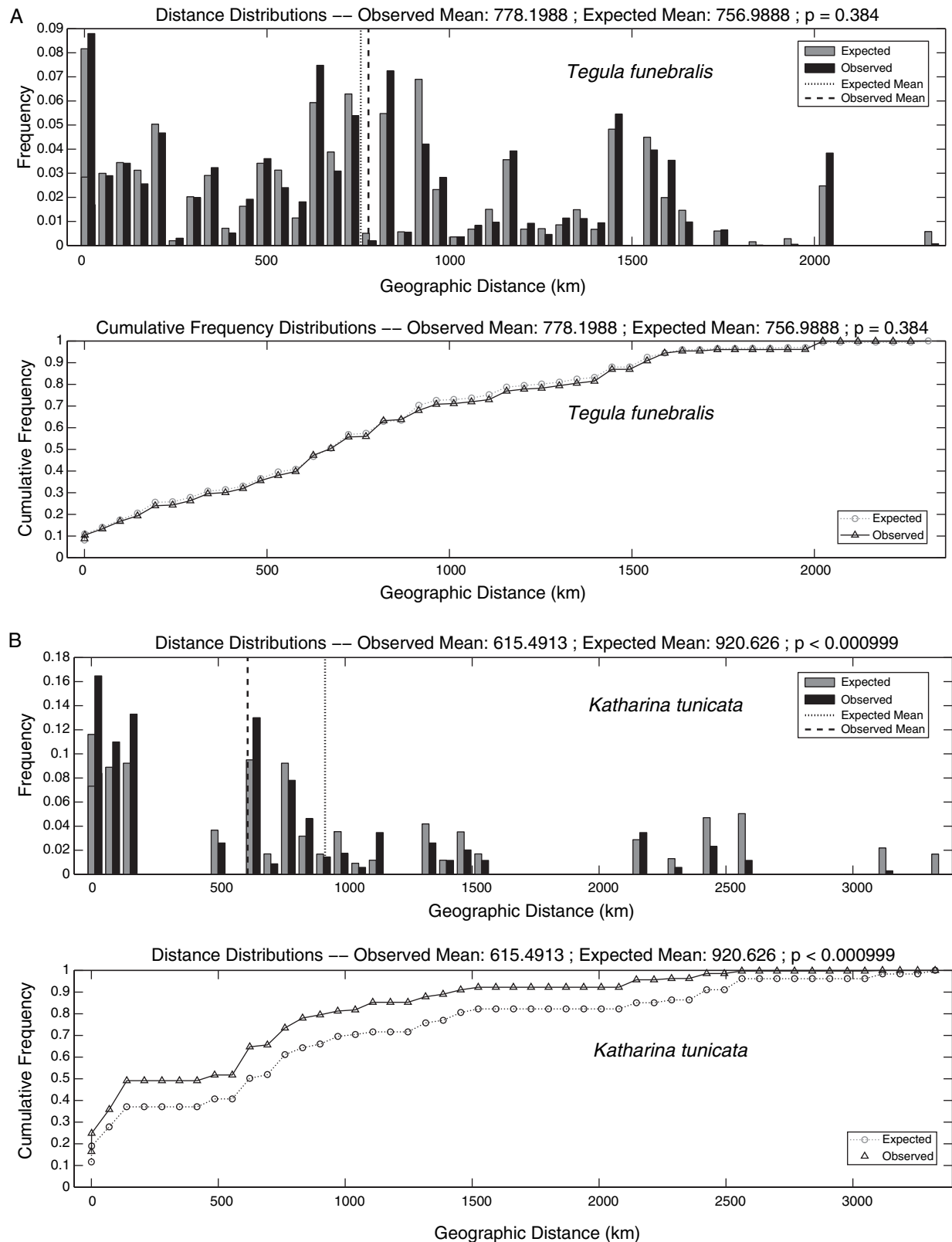
**Figure 4.** Analysis of *Tegula funebralis* (**A**) and *Katharina tunicata* (**B**) COI mtDNA data sets using SAShA. The 2 species have similar within-species patterns of variation by traditional population genetics analysis, however, SAShA reveals that *Katharina*'s haplotypes are significantly underdistributed.

**Table 2** SAShA jackknife results for *Tegula* and *Katharina*

| Allele jacknifed | % of data | OM | EM | P | % Change in OM |
|---|---|---|---|---|---|
| *Tegula funebralis* | | | | | |
| 0 | — | 778.199 | 756.99 | 0.389 | 0 |
| 1 | 32.67 | 744.73 | 756.99 | 0.733 | −4.3008 |
| 2 | 8.67 | 775.162 | 756.99 | 0.488 | −0.3902 |
| 3 | 8 | 782.284 | 756.99 | 0.333 | 0.525 |
| 4 | 4.67 | 780.053 | 756.99 | 0.358 | 0.2382 |
| 5 | 4 | 778.656 | 756.99 | 0.389 | 0.0587 |
| 6 | 2.67 | 778.699 | 756.99 | 0.414 | 0.0643 |
| 7 | 2.67 | 778.906 | 756.99 | 0.371 | 0.0909 |
| 8 | 2.33 | 777.9 | 756.99 | 0.405 | −0.0384 |
| 9 | 1.33 | 777.994 | 756.99 | 0.416 | −0.0263 |
| 10 | 1.33 | 778.34 | 756.99 | 0.411 | 0.0181 |
| 11 | 1.33 | 778.589 | 756.99 | 0.395 | 0.0501 |
| 12 | 1 | 778.339 | 756.99 | 0.413 | 0.018 |
| 13 | 1 | 778.556 | 756.99 | 0.384 | 0.0459 |
| 14 | 1 | 778.056 | 756.99 | 0.413 | −0.0183 |
| 15 | 1 | 778.285 | 756.99 | 0.404 | 0.0111 |
| 16 | 0.67 | 778.046 | 756.99 | 0.419 | −0.0197 |
| 17 | 0.67 | 778.299 | 756.99 | 0.374 | 0.0129 |
| 18 | 0.67 | 778.337 | 756.99 | 0.431 | 0.0178 |
| 19 | 0.67 | 778.173 | 756.99 | 0.38 | −0.0033 |
| 20 | 0.67 | 778.173 | 756.99 | 0.389 | −0.0034 |
| 21 | 0.67 | 778.202 | 756.99 | 0.382 | 0.0005 |
| 22 | 0.67 | 778.093 | 756.99 | 0.411 | −0.0137 |
| 23 | 0.67 | 778.299 | 756.99 | 0.385 | 0.0128 |
| 24 | 0.67 | 778.264 | 756.99 | 0.393 | 0.0084 |
| 25 | 0.67 | 778.337 | 756.99 | 0.402 | 0.0178 |
| *Katharina tunicata* | | | | | |
| 0 | — | 615.491 | 920.63 | 0.004 | 0 |
| 1 | 14.42 | 686.295 | 920.63 | 0.033 | 11.5035 |
| 2 | 13.46 | 581.286 | 920.63 | 0.001 | −5.5574 |
| 3 | 11.54 | **716.354** | **920.63** | **0.12** | **16.3873** |
| 4 | 8.65 | 526.19 | 920.63 | 0 | −14.5089 |
| 5 | 6.73 | 587.68 | 920.63 | 0.011 | −4.5186 |
| 6 | 5.77 | 615.471 | 920.63 | 0.001 | −0.0033 |
| 7 | 3.85 | 615.782 | 920.63 | 0 | 0.0473 |
| 8 | 2.88 | 620.076 | 920.63 | 0 | 0.7448 |
| 9 | 1.92 | 608.09 | 920.63 | 0 | −1.2025 |
| 10 | 1.92 | 616.687 | 920.63 | 0 | 0.1943 |
| 11 | 1.92 | 617.275 | 920.63 | 0.003 | 0.2899 |

Each row represents the results of an analysis done by removing haplotypes sequentially. Note that *Tegula* is never significantly different from panmixia, whereas *Katharina* remains significant after removing any haplotype but the third. This third haplotype is the most important driver of the overall result in the *Katharina* data set, as its removal yields the greatest percent change in OM.

diversity, its performance may suffer relative to $\Phi_{ST}$ with smaller deme sizes. However, genetic subdivision is likely easier to detect at smaller population sizes, so we do not anticipate that this difference between estimators is appreciable in practice. Demes that are small due to a recent bottleneck (and not yet equilibrated, such as might be expected in a conservation context) are likely to be disproportionately diverse. Because SAShA excels with small, highly diverse demes, its relative advantage would be at a peak in that case.

In sum, $\Phi_{ST}$ is an extremely effective statistic in studies with well-sampled, site-based collection designs, but it is much less sensitive with lower density landscape sampling.

Conversely, AAIA shines with sparse landscape sampling when gene flow is high but fails in the more thorough location-based sampling typical of traditional population genetics. SAShA's OM statistic retains the landscape benefits of AAIA while addressing that statistic's high false-negative rate, making SAShA effective for the vast majority of real-world sampling schemes.

### Real-World Data Sets

The 2 real-world data sets, *Tegula* and *Katharina*, illustrate the practical utility of the SAShA approach. These 2 marine species are similar in terms of life history, and the overall within-species patterns of genetic variation appear similar in the 2 using traditional genetic analysis. Yet the haplotypes of *Katharina* are significantly underdistributed, suggesting a limitation to gene flow that was not apparent otherwise. Compared with *Tegula*, *Katharina* may experience greater drift and local selection in the face of considerable gene flow over space. Recent work has shown that, contrary to a common assumption, species with shared reproductive and developmental characteristics can have quite different within-species population genetics patterns (Marko 2004; Kelly and Eernisse 2007). As a result, the differences SAShA reveals between the 2 mollusk species may be generally illustrative of hidden variation in spatial genetic patterns among superficially similar species.

### Strengths and Weaknesses of SAShA

SAShA is more likely than $\Phi_{ST}$ to detect the low levels of genetic subdivision that can exist when gene flow is relatively high (Nm ≈ 100 when $m = 10^{-3}$ and $N = 10^5$). This is a result of analyzing the distribution of individual alleles rather than population-level averages and of simultaneously incorporating geographic and genetic data into a combined analysis. SAShA does not require user-defined populations, making it robust to a wide variety of sampling schemes and eliminating the uncertainty associated with estimating the extent of natural populations, required for most other analyses.

The ability to analyze alleles individually is another strength of our approach. A significant $\Phi_{ST}$ may be due to multiple scenarios (e.g., large pairwise differences between alleles in the data set, nonrandom distribution of similar haplotypes, or a very large number of alleles), whereas the allele-by-allele SAShA makes clear which alleles drive the overall statistics, as well as the amount of discord among alleles in the data set. This method can reveal nonrandom distributions of rarer alleles, which may be more useful indicators of recent gene flow patterns than common, widespread alleles (Slatkin 1985). Using the allele-by-allele analysis, a researcher may discover individual alleles that behave differently from the rest of the data set, perhaps focusing on these for further scrutiny (e.g., to test for selection or for further collection efforts).

Whereas $F_{ST}$ and other statistics do not differentiate metapopulation dynamics from stepping-stone migration or rare long-distance dispersal events, an allele-by-allele SAShA

may aid in distinguishing between models of gene flow. Consistently-underdistributed alleles in a data set suggest a stepping-stone model of gene flow, in which alleles arise via mutation and spread gradually with limited migration. A data set with consistently overdistributed alleles suggests either routine long-distance dispersal among stable populations or a metapopulation in which patches may be colonized from distant source populations. Data sets containing a mixture of random, over-, and underdistributed alleles may be indicative of more complex population genetics or biogeographic histories.

A final advantage of our method is the ability to analyze a variety of data types. Because SAShA makes use of shared alleles among populations, the input data can come from any number of sources (e.g., DNA sequence data, microsatellites, random amplification of polymorphic DNA, restriction fragment length polymorphism, single nucleotide polymorphism [SNPs]). Multilocus data can be accommodated simply by concatenating vertically the input table of haplotypes-by-populations for each locus. However, the test assumes that each row in the input table (i.e., haplotype, allele, etc.) is independent; as such, phase must be known in order to accommodate diploid heterozygote data. Microsatellites, because of their high mutation rate, may be subject to size homoplasy (Estoup et al. 1995), and should therefore be used with caution (see below). Similarly, SNP data with only 2 alleles may violate SAShA's underlying assumption that identical alleles are identical-by-descent when the mutation rate is high: the same allele may arise multiple times. Therefore, SAShA is probably appropriate for SNP data when the mutation rate is low or when the SNPs can be mapped confidently onto a phylogenetic tree and homoplasy avoided.

Our approach is model-free, but is most appropriate when 2 conditions are met: 1) that alleles identical in state are identical by descent and 2) that migration occurs much faster than the combined effects of mutation and drift. Violating either of these conditions is likely to erode SAShA's effectiveness. Applying the analysis to data sets that have high levels of homoplasy, in which alleles are identical in state but not by descent, will produce unpredictable results. If homoplastic alleles are located in the geographic vicinity of one another, there will appear to be more geographic structure in the data set than is actually present; conversely, geographically distant homoplastic alleles will make structure appear artificially low.

The probability of identical alleles arising independently in DNA sequence fragments of any length is exceedingly low, and decreases geometrically as fragment length increases, making homoplasy unlikely to be a problem for the use of sequence data with SAShA. For microsatellite data, high migration rates and recent coalescent times among demes reduce the effect of size homoplasy on the estimation of population differentiation (Rousset 1996; Estoup et al. 2002), thereby minimizing the potential for homoplastic alleles to affect SAShA when geographic genetic differentiation is subtle. If a particular microsatellite locus is known to have a high level of size homoplasy,

however, it may lead to misleading results in any spatial genetics analysis including SAShA. Finally, because any mutational events leading to homoplasy are expected to occur independent of geography, they are not likely to bias the results of the method.

Species with very low gene flow, in which the second condition is violated and effects of drift and mutation outweigh those of migration, may appear to have slightly less geographic structure than they actually do. In practice, data sets with very low gene flow will be obviously structured; it is simply the calculation of the spatial scale over which alleles are shared that will be slightly affected by these local mutants. SAShA's statistics therefore represent a conservative assessment of the amount of geographic structure among homologous alleles.

## Conservation Implications

SAShA's results indicate the scales over which genetic information is shared. Many conservation applications, such as the design of marine reserves, require precisely this kind of information in order to efficiently size and space reserves for maximal protection of genetic diversity. For example, if a continuously distributed species has an OM distance of 500 km between pairs of shared alleles, reserves spaced 200 km apart are likely to be more than sufficient for protecting many of that particular species' alleles, whereas reserves 2000 km apart would be insufficient, probably failing to encompass even some common alleles. Reserve spacing and conservation goals can be quickly evaluated in this way by looking at the geographic distribution of species' individual alleles. Knowing how these alleles are distributed in space also makes possible more sophisticated probabilistic analyses; one could, for example, determine the likelihood of encompassing 95% of a species' alleles given reserves spaced at 300 km. Time series data for a species, such as those being collected annually for *Balanus glandula* along the Oregon coast by the Palumbi laboratory, would make such a calculation much more dynamic, incorporating temporal stochasticity along with the geographic patchiness inherent in population genetics.

Finally, because SAShA is more sensitive than existing statistical methods in many high migration rate species, it will be advantageous in identifying subtle, but real, distinct population segments under the US Endangered Species Act (ESA). Presently, many such claims are undermined by insufficient statistical power: type II error regularly goes undiagnosed in the scientific literature and in turn the relevant federal agencies use the lack of evidence for species' genetic subdivision as evidence of its absence (see, e.g., Brosi and Biber 2009). SAShA is a partial remedy to this problem, if only because it is more likely to detect subtle genetic structure in many real-world sampling scenarios.

To be a discrete population segment (DPS), worthy of protection under the ESA, a biological entity must be 1) discrete, 2) significant, and 3) endangered/threatened. SAShA provides another tool for evaluating the degree to and/or significant. Whereas "discrete" is a more biological

concept, "significant" is a more policy-driven concept—looking at the underlying distribution of alleles gives interested parties more information on which to make their decision as to whether an entity is a DPS. The National Marine Fisheries Service and the U.S. Fish and Wildlife Service, the 2 federal agencies charged with applying the ESA, most often use genetic differentiation at neutral loci to inform the "discreteness" prong of the DPS test. Such differentiation can also imply local adaptive divergence elsewhere in the genome, potentially speaking to the "significance" of a proposed DPS as well, because genetic adaptations contribute to the "evolutionary legacy" of a species (see Waples 2006). SAShA is useful in identifying particular alleles that may be selected for or against, in addition to looking at overall patterns of spatial structure.

## Conclusion

Slight genetic structuring in wild species can be ecologically important but often cannot be detected by conventional methods. Identifying cases of subtle genetic differentiation in the face of gene flow is particularly crucial in species needing management or conservation plans (Waples 1998; Palumbi 2003). We have shown that SAShA is useful for detecting structure in species with relatively high migration rates and small sample sizes, complementing traditional approaches to population genetic analyses. The method is designed to account for diverse sampling schemes, can incorporate a variety of data types, and returns results that are easily interpreted in a geographically explicit context. Finally, analysis of individual alleles provides for a nuanced understanding of the processes underlying the trends observed in the data set. The SAShA MATLAB source code and a downloadable Windows executable program are available at http://sasha.stanford.edu.

## Appendix 1

### The SAShA Algorithm

The algorithm for the calculation of OM and their significance values is as follows:

(1) First, render (m) alleles from (n) locations into an (m × n) allele by location matrix (H), the elements of which represent the number (h) of a given allele (i) in a given location (j):

$$H = \begin{bmatrix} h_{11} & h_{12} & \dots & h_{1n} \\ h_{21} & h_{22} & \dots & h_{2n} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ h_{m1} & h_{m2} & \dots & h_{mn} \end{bmatrix}.$$

(2) Produce an (n × n) pair-wise geographic distance matrix (G) representing the geographic distance (g) between each pair of populations (i and j). All diagonal entries, where i = j, will be zero, and the matrix will be symmetrical along the diagonal ($g_{ji} = g_{ij}$).

$$G = \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1n} \\ g_{21} & g_{22} & \dots & g_{2n} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ g_{n1} & h_{n2} & \dots & h_{nn} \end{bmatrix}.$$

(3) Produce an *observed distance distribution*, (ODD) that consists of a list of geographic distances between all pairs of occurrences of each allele, including those with zero distance (i.e. pairs of the same allele within a given population).
For the $k^{th}$ allele where k = 1 . . . m, and
the $i^{th}$ and $j^{th}$ populations, where i = 1 . . . n and j = i . . . n
  if i ≠ j (i.e., the pair of allele occurrences does not occur in the same population)
    ODD contains H(k,i) * H(k,j) instances of G(i,j)
  if i = j (i.e., the pair of allele occurrences occurs within a population)
    ODD contains (H(k,i) * (H(k,i)-1)) / 2 instances of G(i,i)

(4) Produce an *expected distance distribution*, (EDD) that consists of a list of geographic distances between all possible pairs of samples in the dataset regardless of allele. This serves as our null expectation under random migration.
For the $k^{th}$ allele where k = 1 . . . m, and
the $i^{th}$ and $j^{th}$ populations, where i = 1 . . . n and j = i . . . n
  if i ≠ j (i.e., the pair of allele occurrences does not occur in the same population)
    EDD contains $\Sigma_k$ H(k,i) * $\Sigma_k$ H(k,j) instances of G(i,j)
  if i = j (i.e., the pair of allele occurrences occurs within a population)
    EDD contains ($\Sigma_k$ H(k,i) * ($\Sigma_k$ H(k,i)-1)) / 2 instances of G(i,i)

(5) From these two distributions (ODD & EDD), one can calculate both the observed and expected mean distances (OM and EM, respectively) between shared alleles.
  To calculate the means of the observed distance distribution and the expected distance distribution, one simply takes the arithmetic mean of each distribution. The difference between these means (OM-EM) expresses the geographic distance by which the alleles in the dataset are over- or underdistributed. Overdistributed datasets result in a difference greater than zero; for underdistributed datasets the difference is less than zero. The limits of the difference are determined by the dataset; the statistic may vary between the positive and negative values of the largest geographic distance in a given dataset.

(6) To assess the significance of both the differences between OM and EM, randomly permute the observed allele by location matrix (H) $N_p$ times, maintaining row and column sums constant (for the simulations presented above, $N_p$ = 1000; larger numbers of permutations result in more precise estimates of significance). Recalculate both statistics

(OM-EM) for each of the $N_p$ permuted datasets. Note that because row and column sums are held constant, the null expectation for each permuted dataset remains unchanged. After $N_p$ permutations, compare the original observed statistics to their respective distributions generated by permutation. The proportion of permuted statistics more extreme than the observed statistic serves as the p-value for a significance test.

The above calculations yield a difference statistic (OM-EM) and its respective significance value for the overall dataset. In addition, we extend the analysis in two ways: (a) by applying the algorithm to each allele individually and (b) by jackknifing—repeatedly assessing the overall dataset while sequentially removing each allele in turn. These subsequent analyses provide greater detail and make clear which alleles drive the geographic pattern in the overall dataset. These extensions are calculated as follows:

(7) To analyze each allele in the dataset independently (e.g., the haplotype-by-haplotype analysis presented in Figure A1), carry out steps 3–6 as above, using only the focal allele to calculate ODD (step 3) but retaining the entire dataset to calculate the EDD (step 4). The way in which the statistics (OM-EM) are calculated remains unchanged (step 5), however their significance is calculated based on the statistical distribution calculated for only the focal haplotype out of each permuted dataset (step 6).

(8) To analyze the relative importance of each allele by removing one at a time from the dataset (jackknife method), first apply the algorithm to the overall dataset as above (steps 1–6) to yield OM-EM and its associated significance value.
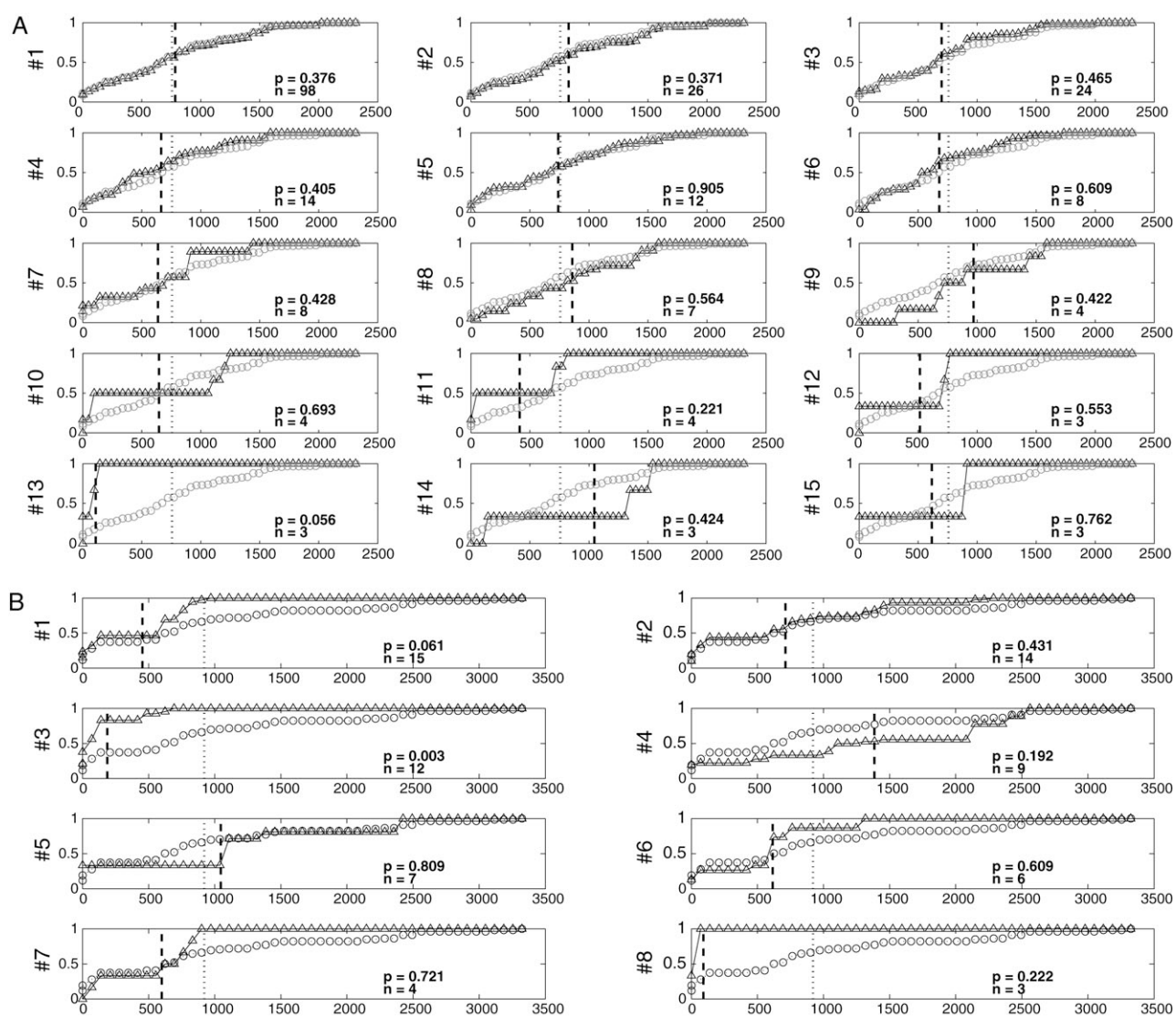


**Figure A1.** Haplotype-by-haplotype analysis of *Tegula funebralis* (**A**) and *Katharina tunicata* (**B**) data sets. The number of occurrences of each haplotype is given (n) as well as the haplotype-specific SAShA statistics. Haplotype number, in order of decreasing frequency, is on the y-axis for each plot. The expected distributions are represented by gray open circles, and observed distributions by dark open triangles.

Then, remove one allele from the input data matrix, H, to create H′. Perform all calculations on H′, and calculate the percentage change for OM-EM for the H′ data matrix relative to that for H. Repeat this procedure for each allele.

## Funding

## Acknowledgments

## References

Beerli P. 2006. Comparison of Bayesian and maximum-likelihood inference of population genetic parameters. Bioinformatics. 22:341–345.

Brosi BJ, Biber EG. 2009. Statistical inference, Type II error, and decision making under the US Endangered Species Act. Front Ecol Environ. 7:487–494.

Epperson BK. 2003. Geographical genetics. Princeton (NJ): Princeton University Press.

Estoup A, Jarne P, Cornuet J-M. 2002. Homoplasy and mutation model at microsatellite loci and their consequences for population genetics analysis. Mol Ecol. 11:1591–1604.

Estoup A, Tailliez C, Cornuet JM, Solignac M. 1995. Size homoplasy and mutational processes of interrupted microsatellites in 2 bee species, *Apis mellifera* and *Bombus terrestris* (Apidae). Mol Biol Evol. 12:1074–1084.

Excoffier L. 2000. Arlequin users manual. [Internet]. Available from: http://anthro.unige.ch/arlequin.

Excoffier L, Laval G, Schneider S. 2005. Arlequin ver. 3.0: an integrated software package for population genetics data analysis. Evol Bioinform Online. 1:47–50.

Excoffier L, Novembre J, Schneider S. 2000. SIMCOAL: a general coalescent program for simulation of molecular data in interconnected populations with arbitrary demography. J Heredity. 91:506–509.

Excoffier L, Smouse P, Quattro J. 1992. Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data. Genetics. 131:479–491.

Folmer O, Black M, Hoeh W, Lutz R, Vrijenhoek R. 1994. DNA primers for amplification of mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates. Mol Mar Biol Biotechnol. 3:294–299.

Hardy OJ, Vekemans X. 1999. Isolation by distance in a continuous population: reconciliation between spatial autocorrelation analysis and population genetics models. Heredity. 83:145–154.

Heywood JS. 1991. Spatial analysis of genetic variation in plant populations. Ann Rev Ecol Syst. 22:335–355.

Hudson RR. 2000. A new statistic for detecting genetic differentiation. Genetics. 155:2011–2014.

Hudson RR, Boos DD, Kaplan NL. 1992. A statistical test for detecting geographic subdivision. Mol Biol Evol. 9:138–151.

Kelly RP, Eernisse DJ. 2007. Southern hospitality: a latitudinal gradient in gene flow in the marine environment. Evolution. 61:700–707.

Manel S, Schwartz MK, Luikart G, Taberlet P. 2003. Landscape genetics: combining landscape ecology and population genetics. Trends Ecol Evol. 18:189–197.

Marko PB. 2004. 'What's larvae got to do with it?' Disparate patterns of post-glacial population structure in two benthis marine gastropods with identical dispersal potential. Mol Ecol. 13:597–611.

Miller MP. 2005. Alleles In Space (AIS): computer software for the joint analysis of interindividual spatial and genetic information. J Hered. 96:722–724.

Moran AL. 1997. Spawning and larval development of the black turban snail *Tegula funebralis* (Prosobranchia: Trochidae). Mar Biol. 128:107–114.

Nei M. 1973. Analysis of gene diversity in subdivided populations. Proc Natl Acad Sci U S A. 70:3321–3323.

Novembre J, Slatkin M. 2009. Likelihood-based inference in isolation-by-distance models using the spatial distribution of low-frequency alleles. Evolution. 63:2914–2925.

Palumbi SR. 2003. Population genetics, demographic connectivity, and the design of marine reserves. Ecol Appl. 13:S146–S158.

Panchal M, Beaumont MA. 2007. The automation and evaluation of nested clade phylogeographic analysis. Evolution. 61:1466–1480.

Peakall R, Smouse P. 2006. GENALEX 6: genetic analysis in Excel. Population genetic software for teaching and research. Mol Ecol Notes. 6:288–295.

Petit RJ. 2007. The coup de grâce for the nested clade phylogeographic analysis? Mol Ecol.

Raymond M, Rousset F. 1995. An exact test for population differentiation. Evolution. 49:1280–1283.

Rousset F. 1996. Equilibrium values of measures of population subdivision for stepwise mutation models. Genetics. 142:1357–1362.

Slatkin M. 1981. Estimating levels of gene flow in natural populations. Genetics. 95:503–523.

Slatkin M. 1985. Rare alleles as indicators of gene flow. Evolution. 39:53–65.

Slatkin M. 1995. A measure of population subdivision based on microsatellite allele frequencies. Genetics. 139:457–462.

Smouse PE, Long JC, Sokal RR. 1986. Multiple regression and correlation extensions of the Mantel test of matrix correspondence. Syst Zool. 35:627–632.

Smouse PE, Peakall R. 1999. Spatial autocorrelation analysis of individual multiallele and multilocus genetic structure. Heredity. 82:561–573.

Strathmann MF. 1987. Reproduction and development of marine invertebrates of the Northern Pacific Coast: data and methods for the study of eggs, embryos, and larvae. Seattle (WA): University of Washington Press.

Templeton AR. 1998. Nested clade analyses of phylogeographic data: testing hypotheses about gene flow and population history. Mol Ecol. 7:381–397.

Templeton AR. 2004. Statistical phylogeography: methods of evaluating and minimizing inference errors. Mol Ecol. 12:789–809.

Templeton AR. 2008. Nested clade analysis: an extensively validated method for strong phylogeographic inference. Mol Ecol. 17:1877–1880.

Vekemans X, Hardy OJ. 2004. New insights from fine-scale spatial genetic structure analyses in plant populations. Mol Ecol. 13:921–935.

Waples RS. 1998. Separating the wheat from the chaff: patterns of genetic differentiation in high gene flow species. J Hered. 89:438–450.

Waples RS. 2006. Distinct population segments. In: Scott JM, Goble DD, Davis FW, editors. The Endangered Species Act at thirty: conserving

biodiversity in human-dominated Landscapes. Washington (DC): Island Press. p. 127–149.

Weir BS, Cockerham CC. 984. Estimating F-statistics for the analysis of population structure. Evolution. 38:1358–1370.

Wright S. 1951. The genetical structure of populations. Anna Eugen. 15:323–354.

Wright S. 1965. The interpretation of population structure by *F*-statistics with special regards to systems of mating. Evolution. 19:395–420.