21 **Title**: Detecting spatial dynamics of range expansions with geo-referenced genome-wide SNP

22 data and the geographic spectrum of shared alleles

23 **ABSTRACT**

24     Uncovering the spatial dynamics of range expansions is a major goal in studies of

25 historical demographic inference, with applications ranging from understanding the evolutionary

26 origins of domesticated crops, epidemiology, invasive species, and understanding species-level

27 responses to climate change. Following the surge in advances that make explicit use of the

28 spatial distribution of genetic data from geo-referenced SNP variants, we present a novel

29 summary statistic vector, the geographic spectrum of shared alleles (GSSA). Using simulations

30 of two-dimensional serial expansion, we find that the information from the GSSA, summarized

31 with Harpending's Raggedness Index (RI), can accurately detect the spatial origins of a range

32 expansion under serial founder models, even with sparse sampling of only ten individuals. When

33 applying to SNP data from two species of the holarctic butterfly genus *Lycaeides*, the suggested

34 origins of expansion are consistent with hindcasts obtained from ecological niche models

35 (ENMs). These results demonstrate the GSSA to be a useful exploratory tool for generating

36 hypotheses of range expansion with genome-wide SNP data. Our simulation experiments suggest

37 high performance even with sampling found in studies of non-model organisms (one sampled

38 individual per location, no outgroup information, and only 5,000 SNP loci).

39 **Keywords:** demographic history, range expansion, geographic spectrum of shared alleles,

40 summary statistic

## Introduction

41　Geographic range shifts are a prevalent feature of virtually all species' histories. They

42　often follow or accompany the origin of species (Gaston 2003), and frequently happen in

43　response to changes in climate, landscape, or biotic opportunities (Konečný *et al.* 2013; Roberts

44　& Hamann 2015). Despite range expansions being a pervasive feature across many species

45　histories, the net outcomes of the interactions between species-specific traits and changes in

46　suitable habitats are poorly understood with regards to changes in the geographic distribution of

47　genetic diversity, adaptability, and community structure (Pharo & Zartman 2007; Thuiller *et al.*

48　2008; Phillips *et al.* 2010). For example, although serial-range expansions are likely a common

49　feature of species histories in the context of past global changes (Excoffier *et al.* 2009), there is

50　still uncertainty about how they affect differentiation and diversification, or how they influence

51　the risk of extinction under future environmental shifts (Charles & Dukes 2008; Fordham *et al.*

52　2014). Tools that uncover the spatio-temporal dynamics of species' ranges are therefore critical

53　for understanding how geographic range dynamics affect the spatial distribution of genetic

54　diversity through drift and selection (Slatkin 1987; Kirkpatrick & Barton 1997; Peischl *et al.*

55　2013). They are also central to forecast species' responses to ongoing and future environmental

56　changes (Brown *et al.* 2016).

57　Among recent attempts to couple genetic and spatial information to infer historical range

58　dynamics, Ramachadran *et al.* (2005) proposed to identify the geographic origin of an expansion

59　using heterozygosity clines that result from the reduction of heterozygosity associated with the

60　sequential founder events that are characteristic of range expansions (Ramachandran *et al.* 2005;

61　Li *et al.* 2008; DeGiorgio *et al.* 2009; Henn *et al.* 2012). Based on the premise that

63    heterozygosity should continuously decrease with distance from the expansion source, this

64    method fits a linear regression between expected heterozygosity at multiple sampled locations

65    and their geographic distance from all possible candidate sources. While this is simple to

66    estimate, this approach has been found to perform well under a relatively narrow range of

67    circumstances, with its accuracy decreasing as the speed of the expansion decreases, or as the

68    time since the end of the expansion increases (Peter & Slatkin 2015). Previously, researchers

69    have used principal component analysis (PCA) of geo-referenced samples to identify the

70    direction of expansion under the assumption that the first principal component—which accounts

71    for the greatest amount of variation—aligns with the expansion source-front axis (Menozzi *et al.*

72    1978). However, this assumption has been challenged by recent evidence indicating that the

73    geographic alignment of the first component may be driven by the geographic distribution of

74    samples, allele surfing, or mathematical artifacts (Novembre & Stephens 2008; François *et al.*

75    2010; Pemberton *et al.* 2013).

76         More recently, Peter and Slatkin (2013) proposed a new statistic, the directionality index

77    ($\psi$), to detect the origins of range expansions using joint geo-referenced and genome-wide

78    information. This statistic takes advantage of allele frequency changes driven by the sequential

79    founder events that often characterize expansions. Specifically, $\psi$ detects asymmetries in 2D site

80    frequency spectra of shared neutral non-ancestral alleles between pairs of populations, under the

81    theoretical expectation that populations further away from the expansion have experienced more

82    genetic drift than populations near the expansion source (Slatkin & Excoffier 2012; Peter &

83    Slatkin 2013). This metric has been also incorporated into an ABC framework, whereby He et al.

84    (2017) used it to estimate the source of range expansion in the Collared pika in Alaska.

4

85    While analyses using the directionality index (Peter & Slatkin 2015) have been highly

86    promising, the relatively high number of sampled individuals, and the requirement of polarized

87    SNPs, make its implementation difficult in non-model organisms, which are typically

88    characterized by limited numbers of geo-referenced genotypes, or for which a reliable outgroup

89    for polarizing SNPs has not been identified. Here we develop a summary statistic vector, which

90    we call the geographic spectrum of shared alleles (GSSA, hereafter), to harvest information from

91    the spatial distribution of shared allelic variants to help identify the geographical origin of range

92    expansions. Our approach is targeted for non-model organisms as it does not require information

93    about ancestral states, requires a single individual per sampled location, and can accommodate a

94    reasonable number of SNPs.

95    We formally describe the GSSA summary statistic vector and use two-dimensional

96    simulations of different expansion histories to statistically explore its behavior. We also

97    demonstrate its applicability with genome-wide SNP data sampled from a set of taxa of the

98    holartic butterfly genus *Lycaeides* distributed in Western North America (Gompert *et al.* 2010).

99    Our analyses demonstrate that the GSSA can be a powerful exploratory tool to identify the

100   relative position of samples along an expansion axis. In *Lycaeides*, the GSSA places the origin of

101   expansions in geographical areas that are consistent with species distribution models for

102   hypothesized ranges shifts during the last glacial maximum.


103   **Materials and Methods**

104   **The geographic spectrum of shared alleles (GSSA)**

105    The Geographic Spectrum of Shared Alleles (GSSA) is a property of—and calculated

106    for—each sampled locality. It captures the geographic distribution of shared co-ancestry with

107    other localities by summarizing the geographic distances between copies of the minor allele

108    observed at each SNP site (that is, the allele with the smallest frequency in the global sample for

109    each site; Nielsen et al. 2012). As such, the GSSA of a given location is a histogram that depicts

110    the cumulative frequency distribution of the geographic distances separating each minor allele

111    shared between the focal location and all other sampled locations across all SNPs (Fig. 1). On an

112    intuitive level, this new summary statistic vector can be understood as a spatial structure function

113    that captures the relative change in allelic similarity from a focal location as a function of

114    geographic distance. Importantly, the GSSA is a location-specific statistic which differs from

115    other structure functions such as spatial correlograms that are globally calculated. These latter

116    statistics quantify the spatial dependence of an observed variable across the entire set of

117    observations using a covariance matrix partitioned into distance classes (Smouse & Peakall

118    1999; Legendre & Legendre 2012). Instead, each location-specific GSSA is directly based on

119    distances between minor alleles shared with a focal locality.

120    The information captured in the GSSA helps identify the relative position of each

121    sampling location with respect to the expansion source-front axis, because the similarity in

122    genetic constitution between locations depends on the extent of shared history under an

123    expansion history (Peter & Slatkin 2015; Bradburd *et al.* 2016). Because the amount of genetic

124    drift increases with distance to the source of an expansion (Slatkin & Excoffier 2012; Peter &

125    Slatkin 2013), and because colonization through different geographic paths should lead to allelic

126    segregation (Hallatschek *et al.* 2007; Knowles & Alvarado-Serrano 2010; François *et al.* 2010),

127    the amount of shared genetic variants can potentially further inform about the direction of the

128    expansion. As such, if this information is captured by our summary statistic vector, it could be

129    useful in a simulation-based statistical method—such as approximate Bayesian computation or

130    supervised machine learning—for testing alternative expansion histories (Pudlo *et al.* 2015;

131    Joseph *et al.* 2016; Schrider & Kern 2018).


132    **GSSA overview**

133    The calculation of the set of GSSAs, one per unique sampled geo-referenced location,

134    requires two sets of data: i) the geographic distances between all sampled locations, and ii) a

135    record of the counts of minor allele copies of each SNP locus at each sampled location (Fig. 1a).

136    Here, we restrict our discussion to geographic distance between samples, but note that users

137    could use an alternative distance metric, such as river-distance, when appropriate. Below we

138    verbally summarized the steps involved in the calculation of the GSSA for one sampled location

139    (a formal mathematical description is presented in the next section).

140    First, using the entire set of geographic distances between all sampled locations, the

141    Sturges (1926) equation is deployed to identify the optimal number of geographic distance

142    classes and breakpoints for the construction of all of the location-specific GSSAs. This Sturges

143    binning scheme is then applied to the set of geographic distances associated with the *i*th focal

144    sampled location, such that there will be a specific histogram ($h_{geo_i}$) for each of the locations

145    (Fig. 1e). Second, a matrix of minor allele's presence/absence ($G_i$) is constructed for the *i*th focal

146    sampled location (Fig. 1a). Third, a vector ($\vec{S}_i$) is constructed for the *i*th focal sampled location

147    from the geographic distribution of minor alleles. Specifically, each location-specific $\vec{S}_i$ vector

148    lists the geographic distances between each minor SNP allele present at the focal sampling

149    location and all other copies of the same minor allele in the entire sample (Fig. 1c). Fourth, this

150    location-specific vector is then converted into a corresponding location-specific histogram ($h_{gen_i}$

151    ) using the Sturges' (1927) binning scheme previously defined (Fig. 1d). Finally, to make our

152    statistic independent of the specific sampling scheme (i.e., sampling locations and their influence

153    on the binning scheme), the location-specific histogram constructed from the vector of distances

154    between minor alleles ($h_{gen_i}$) is regressed against the values from the corresponding

155    "geography-only" histogram ($h_{geo_i}$; i.e., the "null" distribution of geographic distances). The

156    residuals from this regression are finally used to build the location-specific histogram for the $i$th

157    focal location (GSSA$_i$).


158    **GSSA calculation**

159    Formally, the two sets of data extracted from the series of georeferenced sampled

160    genotypes can be defined as two matrices. With one individual per sampled location, $I$ locations,

161    $M$ sets of chromosomes (i.e., ploidy), and $L$ loci (i.e., SNPs), the arrangement of minor alleles in

162    the entire sample relative to each location can be characterized by a $M$ x $L$ genotype matrix, $G_i$,

163    where each element $G_{i_{m,l}}$ takes the value of 0 or 1, depending on whether a copy of the minor

164    allele is present (1) or not (0) at each individual's SNP locus and DNA strand (note that phasing

165    is not necessary because the method treats each locus independently; Fig. 1a). In turn, the set of

166    Euclidean or effective (Shirk & Cushman 2014; Davis $et$ $al.$ 2018) geographic distances between

167    all locations can be characterized by an $I$ x $I$ square distance matrix, $D$, where $D_{i,i} = 0$ and $D_{i,j} =$

168    $D_{j,i}$. From this latter matrix, a set of location-specific geographic-distances vectors, $\vec{d_i}$, can be

169     obtained by subsetting the D matrix so that $\vec{d}_i = D_{i,1:I}$ (Fig. 1b). Altogether, these data are used

170     to construct a GSSA summary statistic vector for each sampled location following five steps:

171        *Step 1*. First Sturges' (Sturges 1926) equation is applied to the set of pairwise geographic

172     distances between all localities contained in matrix *D* to identify an optimal series of consecutive

173     distance classes that are later used for histogram binning. We chose Sturges' (1926) equation for

174     its simplicity and common use in biological research (Ramírez-García *et al.* 1998; Fagua &

175     Gonzalez 2007; Pires *et al.* 2016; Cardozo *et al.* 2018). Although alternative binning schemes

176     would in theory be possible, an exploration of a coarser binning scheme (Scott 1979) and the

177     more granular Rice rule (Jones *et al.* 2001--) showed that although Sturge's scheme resulted in

178     relatively better performance with regards to the dynamics of the GSSA and identification of

179     expansion colonization history , accuracy was good across all three binning schemes (Supp. Fig.

180     1).

181        *Step 2*. Each genotype matrix $G_i$ is used to construct a location-specific vector, $\vec{S}_i$, which

182     summarizes the aggregated relative spatial distribution of minor alleles at location *i* (Fig. 1c).

183     The elements of this vector are defined as:

184

$$s_{i_{m,l}} = \begin{cases} \emptyset & when \ \left( G_{i_{m,l}} * G_{j_{m,l}} \right) = 0 \\ D_{i,j} & when \ \left( G_{i_{m,l}} * G_{j_{m,l}} \right) \neq 0 \end{cases}; \quad s_{i_{m,l}} \in \vec{S}_i \ \bigwedge \ j \in \{1:I\}$$

[1a];

185     where $\emptyset$ = null, $G_{i_{m,l}}$ = element of individual *i*'s genotype matrix that denotes the presence or

186     absence of a minor allele at locus *l*, strand *m*, and $D_{i,j}$ = geographic distance between localities *i*

187     and *j*. Elements in the $\vec{S}_i$ vector with a value of zero correspond to instances in which the same

188     minor allele at a SNP locus is homozygous within an individual or present in more than one

189   individual sampled from the same locality because the geographic distance between these allele

190   copies would be zero ($D_{i,i} = 0$). On the other hand, non-zero elements correspond to the

191   geographic distances between individuals from different localities sharing the same minor allele

192   at a SNP locus (Fig. 1c). Effectively, $\vec{S}_i$ can also be derived from the frequency of the minor

193   allele at each locus and location by aggregating over DNA strands so that the set of $\vec{S}_{i_{m_{1:M},l}}$

194   elements of this vector are defined as:

195

$$s_{i_{m_{1:M},l}} = \left\{ \underbrace{\langle D_{i,j} \rangle}_{(f_i \times f_j)\ times} \right\}; \qquad s_{i_{m,l}} \in \vec{S}_i \ \wedge \ j \in \{1:I\}$$

[1b];

196   where $f_i$ is the frequency of copies of the minor allele at location $i$, locus $l$ and $D_{i,j}$ as in eq. [1a].

197       *Step 3.* After removing all null elements (Ø) from vector $\vec{S}_i$ (Fig. 1c), which number

198   depends on the number of loci at which each locality does not contain a minor allele copy (Fig.

199   1a), the size of each $\vec{S}_i$ varies among locations. Therefore, it is necessary to compress each $\vec{S}_i$

200   into a vector with an equal number of elements for all individual sampled locations. To do this,

201   each vector $\vec{S}_i$ (Fig. 1c) is converted into a histogram by calculating the frequency of $\vec{S}_i$

202   elements that fall within each of the distance classes previously determined in *Step* 1. The

203   resulting location-specific geo-genetic histogram ($h_{gen_i}$) summarizes the relative spatial

204   distribution of minor alleles copies across loci present at each ($i^{th}$) sampling location (Fig. 1d).

205       *Step 4.* To correct for the "null" geographic expectation introduced by the specific

206   geographic position of each sample, each location-specific vector of geographic-distances, $\vec{d}_i$, is

207   first converted into a histogram by again estimating the frequency of observations that fall within

208    each of the same set of distance classes previously used to create the geo-genetic histograms.

209    This generates a unique geographic histogram ($h_{geo_i}$) for each ($i^{th}$) location that is analogous to

210    the corresponding geo-genetic histograms ($h_{gen_i}$) (Fig. 1e), yet that only containing geographic

211    information.

212         *Step 5*. Finally, the elements of each geo-genetic histogram ($h_{gen_i}$) are regressed against

213    the corresponding elements of the geographic-distance histogram ($h_{geo_i}$) (Fig. 1f), as defined

214    here:

$$\hat{h}_{gen_{i[b_{k=1}:b_{k=K}]}} = \beta \left( h_{geo_{i[b_{k=1}:b_{k=K}]}} \right) + \varepsilon \quad \text{[2];}$$

216    where $h_{gen_i}$ = geo-genetic histogram for location $i$, $h_{geo_i}$ = geographic histogram for location $i$,

217    [$b_k$] = histogram's distance class, $b$, ranging from k = 1 to k = K (i.e., the maximum number of

218    bins as determined by Sturges' equation), $\beta$ = simple regression coefficient, and $\varepsilon$ = error term.

219         The vector of absolute residuals resulting from each regression (eq. [2]) constitutes the

220    location-specific spatial summary statistic vector we named the GSSA and it is defined for each

221    location $i$ as the absolute difference between the location-specific regression-predicted values

222    (eq. [2]) and the location-specific geo-genetic histogram values:

$$GSSA_i = \left| \hat{h}_{gen_{i[b_{k=1}:b_{k=K}]}} - h_{gen_{i[b_{k=1}:b_{k=K}]}} \right| \quad \text{[3];}$$

224    with $\hat{h}_{gen_i}$, $h_{gen_i}$, and [$b_k$] defined as in eq. [2].

225    **GSSA Implementation**

226        All five steps described above can be implemented from a list of genotypes and a list of

227    individuals' sampling coordinates, or optionally a set of distances between genotypes in case

228    users choose to use effective distances instead (Shirk & Cushman 2014; Davis *et al.* 2018). The

229    aggregated relative spatial distribution of minor alleles vector for each locality, $\vec{S}_i$, can be

230    obtained from the multi-site frequency spectrum (multiSFS), which summarizes the frequency of

231    shared allele variants across populations, and the set of distances among all sampled locations.

232    This multiSFS can be calculated from the list of genotypes in available programs such as **δaδi**,

233    whereby the multiSFS should be polarized based on the minor allele where each sampled

234    geo-referenced location is considered a "population" in the multiSFS (Gutenkunst *et al.* 2009).

235    Scripts to estimate the multiSFS along the rest of necessary steps to calculate the GSSAs for each

236    location are available at https://bitbucket.org/diegofalvarado-s/gssa_v0.0/src. If multiple

237    individuals per location are available, the method can be adapted by aligning all genotypes per

238    location into a single "polyploid individual" so that individuals' genotypes are treated as

239    different DNA strands. Similarly, polyploid individuals are likewise accommodated as the

240    method treats each SNP locus independently and thus, no phasing is necessary.

241    **Raggedness index**

242        A useful property of the GSSA is that its shape can potentially carry information about a

243    location's relative age of colonization and its historical connectivity with other populations given

244    an expansion history. To explore and summarize the overall behavior of the GSSA with respect

245    to expansion sources, we use Harpending's raggedness index (Harpending 1994). Harpending's

246    raggedness index (RI) was originally introduced to quantify the shape of the histogram obtained

247    from the average pairwise genetic differences amongst individuals in the context of inferring the

248    history of population size change due to the predictive relationships between the shape and the

249    timing and magnitude of size change. Here we repurpose it to condense the GSSA elements into

250    a single variable that is correlated with the shape of the GSSA and thus, summarizes the behavior

251    of the GSSA and its spatial and temporal relationship with the expansion source. In our

252    calculation of each locality-specific RI, we disregard all distance classes within each $GSSA_i$ that

253    are equivalent to the distance classes in the corresponding $h_{geo_i}$ that have a value of zero (which

254    arise for distance classes not involving the focal locality). We then correct the RI for the number

255    of comparisons used for its calculation:

256
$$RI = \frac{\sum_{b_k=2}^{b_k=K}(b_k - b_{k-1})^2}{K' - 1}; \qquad b_k \in GSSA_i \mid h_{geo_{b_k}} \neq 0 \qquad [3],$$

257    where $b_k$ indicates a distance class and $K'$ the total number of distance classes used in the

258    calculation.

259        Locations further away from the serial-expansion source are expected to share more

260    genetic variants with nearby locations that were colonized through the same expansion path

261    (Knowles & Alvarado-Serrano 2010), and hence should tend to present a left-skewed and more

262    ragged GSSA (Supp. Fig. 2a). This is because, in this case, short distance bins should be overly

263    represented with a marked drop at intermediate bins, thereby inflating the RI from the

264    "neighborhood effect" resulting from allele surfing (Hallatschek & Nelson 2008; François *et al.*

265    2010). This process can cause genetic variants at the expansion front to increase in frequency

266    within constrained areas of derived populations (Excoffier *et al.* 2008). In contrast, locations

267    closer to the source should tend to have a more uniform and non-skewed GSSA histogram

268    because these locations are expected to retain the shared genetic variants present at different

269    frequencies in other locations in the sample. In this latter scenario (Supp. Fig. 2b), most variants

270    come from standing genetic variation in the source population due to the comparatively smaller

271    amount of drift experienced by these close-to-source populations (Excoffier *et al.* 2008; Peter &

272    Slatkin 2013).

**Simulation experiments**

274    To evaluate the behavior of the GSSA, we simulated 5,000 unlinked SNPs sampled from

275    10 spatially separated individuals under different serial range expansion histories, where 99

276    demes are colonized by a single source deme. Spatially implicit simulations were conducted

277    using fastsimcoal v2.5.2 (Excoffier & Foll 2011; Excoffier *et al.* 2013), with each simulation

278    starting with a source deme at either of four different starting locations (Fig. 2). At each

279    sequential expansion step going forward in time, we set a proportion of individuals within each

280    deme (*f*) to move out from the demes they occupied into an immediate neighboring deme,

281    excluding diagonal movements (i.e., 4-neighbors). We set newly colonized demes to then grow

282    to a size $N_K$ (based on the carrying capacity parameter), in $\tau_r$ generations. After $\tau_c$ generations, we

283    set the next set of demes to be colonized from previously colonized populations, and cycles of

284    serial colonization proceed until the last of 99 demes is colonized. We set the simulations to then

285    run for $\tau$ additional generations after the last colonized deme has reached $N_K$. Through the entire

286    simulation, we allowed colonized demes to exchange individuals with neighboring colonized

287    demes (4-neighbors) according to a migration parameter (*m*) that reflects the per individual per

288    generation probability of migration. After all available demes are colonized, we let the

14

289     simulations continue under a stepping-stone model (Kimura & Weiss 1964) for a number of

290     generations determined by parameter ($\tau$), which signifies the time after initial expansion has

291     ended. Table 1 lists all parameters controlling the simulations.

292          The particular effect of each parameter on the ability of the GSSA vector (summarized

293     with the RI) to locate the expansion source was quantified independently. For that, we kept all

294     parameters—except the one of interest of interest—fixed across simulations (Table 1). For each

295     scenario/parameter combination, we ran 1,000 simulations for a total of 48,000 simulations (1

296     scenario x 4 possible sources x 4 parameters x 3 parameter values x 1,000 simulations). After

297     each simulation was completed (as defined by $\tau$), we gathered a dataset of 5,000 homologous

298     biallelic SNPs from one individual from each of the 10 locations that were randomly selected

299     before the simulations, including the source location. We then used our custom python pipeline

300     to calculate a GSSA (eqs. [2, 3]) and associated RI (eq. [4]) for each of the 10 sampled locations.

301     The source of the expansion was identified based on the distribution of magnitudes of the RI

302     across sampling locations, with the smallest RI value assumed to correspond to the expansion

303     source. Accuracy was measured as the proportion of simulations that correctly identified the

304     source. To capture the empirical reality that the actual expansion source may not have been

305     sampled, we repeated our inferences while excluding the source location from our sample (Fig.

306     S2). In this latter case, accuracy was measured as the proportion of simulations that correctly

307     identified the most proximate location to the source location among those sampled (Supp. Fig.

308     3). In the case of the simulations that include the source, the root mean squared error (RMSE)

309     was also calculated. Additionally, we visualized the magnitude of change in the RI of each

310     sampled location as a function of when each location was colonized.

311     As a point of comparison, we also used Peter and Slatkin's (2013) approach to identify

312     the expansion source on the same simulated datasets. Because Peter and Slatkin's approach make

313     use of interpolation, which reduces the chances of estimating the precise coordinates of the

314     simulated source, we scored a prediction as accurate whenever Peter and Slatkin's directionality

315     index was able to locate the source within the perimeter defined by the 8-neighbor demes around

316     the actual source.

## Empirical Application

318     To illustrate and demonstrate the utility of our approach, we applied it to two taxa of the

319     holarctic butterfly genus *Lycaeides* that are distributed in western North America (Gompert *et al.*

320     2014b). These butterfly taxa are host-specialists and relatively poor dispersers, with distributions

321     that are tightly linked to their host plants (Forister *et al.* 2011; Gompert *et al.* 2014a). Their

322     respective geographical ranges were likely impacted by Pleistocene climatic events (Thompson

323     *et al.* 1993; Thompson & Anderson 2000; Pierce *et al.* 2004), which presumably involved range

324     expansions from refugia since the Last Glacial Maximum (LGM). In particular, we focused on

325     two high-elevation undescribed taxa (Alpine and Jackson *Lycaeides*; Gompert et al. 2014). These

326     data are well-suited for demonstrating our approach because they have a relatively

327     well-characterized evolutionary history and spatial distribution (Gompert *et al.* 2006; Nice *et al.*

328     2013), as well as a thorough, spatially widespread genetic sampling (Gompert et al. 2014) in the

329     western mountains of North America (Alpine *Lycaeides*: 10 localities and 8097 independent

330     homologous SNPs; Jackson *Lycaeides*: 11 localities and 9074 independent homologous SNPs).

16

331    To compare our geographic predictions of source locality with the plausible geographic

332    origin of each taxon's putative range expansion after the LGM, we used ecological niche models

333    (ENMs), projected onto past climates, as independent hypotheses of colonization history under

334    the assumption of Grinnellian niche conservatism. Briefly, we generated an ENM for each taxon

335    in the R package dismo (Hijmans 2012) using a maximum entropy approach (Phillips *et al.*

336    2006) and 19 interpolated climatic surfaces at 30 arc-seconds resolution that summarize global

337    patterns of temperature, precipitation, and seasonality (Hijmans 2012). We only included

338    localities for which genetic data are available (Table 1, Gompert et al. 2014) given the pending

339    taxonomic status of these taxa, which can lead to misidentification issues in available locality

340    databases. To reduce possible sampling bias, for each taxon we filtered the localities based on a

341    minimum distance required between localities. The exact distance used for each species was

342    determined by a variogram approach that establishes the distance over which spatial

343    autocorrelation in environmental conditions is minimal (Brown *et al.* 2016). In addition, to

344    maximize the fit of the model and to avoid unnecessary complexity, we tuned our models using

345    the R package ENMeval (Muscarella *et al.* 2014). The tuning procedure consisted of assessing

346    optimal model parameters based on jackknife cross-validation of the samples (Radosavljevic &

347    Anderson 2014; Muscarella *et al.* 2014). We chose a jackknife approach given the small sample

348    size of the filtered collection points (Galante *et al.* 2018). Model fit was evaluated under different

349    combinations of model features and regularization multipliers (Table S1) using three sequential

350    criteria: omission rate, AUC, and model feature class complexity . Tuning results are

351    summarized in Supp. Table 1. We then hindcasted this ENM to the LGM using climatic CCSM

352    (Community Climate System Model) estimations derived from the Paleoclimate Modelling

17

353   Intercomparison Project Phase II (Braconnot *et al.* 2007). Finally, following Knowles and

354   Alvarado-Serrano (2010), we identified the area(s) of contiguous high predicted suitability (i.e.,

355   refugial distribution) as potential expansion source(s), using as an arbitrary cut-off the top 10[th]

356   percentile of the predicted suitability scores. We then contrasted these hypothesized sources with

357   those inferred by our approach.

## Results

### Simulation experiments

360      Our simulation study showed that the magnitude of the RI calculated from the elements

361   of our GSSA vector was highly correlated with i) the distance from the location of the expansion

362   source, and ii) how long, after the expansion started, a location was colonized (Fig. 3). The

363   metric was able to accurately identify the geographical source of expansion given that one of the

364   sampling localities was the source of expansion, or the first location colonized among those

365   sampled (Fig. S3). Still, the degree of accuracy somewhat depended on the location of the source

366   deme such that lower accuracy occurred when the source was located in a peripheral area (Fig.

367   4). Still, incorrect inferences tended to identify sampling localities proximate to the actual

368   expansion locality, regardless of the location of the expansion source (Fig. 5).

369      By summarizing the GSSA with Harpending's (1994) RI, we were able to correctly

370   identify the source population up to 70% of the time, with RMSE values that approached zero

371   (Fig. 6). However, this accuracy declined with larger numbers of colonists ($f$) and time since all

372   demes have been colonized ($\tau$). Likewise, under an IBD model that arose when enough time had

373   accrued for the signature of the expansion to have significantly diminished (Fig. 6), our metric

18

374     was unable to accurately identify the source population. Our approach did strongly outperform

375     Peter and Slatkin's directionality index ($\psi$) by a margin of up to 60%, except in those

376     aforementioned cases such as a high number of individuals colonizing demes (Fig. 6).

377     **Empirical application**

378         The application of our coupled GSSA-RI approach to the two *Lycaeides* datasets suggests

379     an expansion spatial dynamics that is well in line with species ranges at the LGM according to

380     our ENM hindcasts (Fig. 7). For the Alpine *Lycaeides* distributed in northern Rocky mountain

381     areas, the predicted source coincided with the southernmost samples, which are located in the

382     largest LGM refugium predicted by the ENM hindcasts (Figure 7A). Although the hindcasted

383     range at the LGM climate for this species covers a wide area that encompasses all of the sampled

384     locations, the sampling location predicted by the GSSA coincides with an area of highest

385     suitability. The RI calculated on the GSSA for Jackson *Lycaeides* samples suggested a source

386     also within the ENM hindcast (which was more heterogeneous in comparison with Alpine

387     *Lycaeides* hindcast), and close to one potential refugium (i.e., continuous area of high suitability)

388     (Fig. 7B).

389     # Discussion

390         We introduce a new summary statistic vector, the geographic spectrum of shared alleles

391     (GSSA), which makes joint use of geographic and genetic information. We demonstrate that it

392     can be used to help infer the geographic dynamics of a range expansion. We use Harpending's

393     (1994) RI calculated on the GSSA elements of each sampled genotype to summarize spatial

394    gradients in the distribution of minor alleles that are indicative of the relative position of each

395    genotype along an expansion axis. However, we envision the elements of the GSSA vector may

396    be directly incorporated within a supervised machine learning or ABC inferential framework to

397    obtain parameter estimates associated with source locality likelihoods, and to test competing

398    expansion hypotheses (He *et al.* 2017b; Fraimout *et al.* 2017; Schrider & Kern 2018).

399          However, there are some important assumptions to first examine before deployment of

400    the GSSA. First, as with many other population genomic approaches, it should be used after

401    population structure is explored (Patterson *et al.* 2006; Frichot & François 2015; Petkova *et al.*

402    2016). Specifically, the GSSA is applied to groups of samples suggested to derive from single

403    sources by various other exploratory methods applied cautiously (House & Hahn 2017; Elleouet

404    & Aitken 2018). For example, as late Pleistocene range expansions and recent species invasions

405    can both have multiple origins (Miraldo *et al.* 2011; Ruiz-Cooley *et al.* 2013; Cristescu 2015),

406    care needs to be taken to detect the possibility of admixture from multiple origins *a priori*. In this

407    case, if exploratory methods suggest a species is heavily structured, with a history of expansions

408    from more than one refugium, researchers could apply the GSSA to each sub-sample

409    independently. This may, however, prove a challenging task in certain circumstances because

410    smaller sample sizes (resulting from splitting up the range) may lead to reduced accuracy.

411    Difficulties may also arise if more than two sources are close together geographically, multiple

412    expansions occurred, or if gene flow directions are biased given complex heterogeneous

413    geography (Branco *et al.* 2018; Lundgren & Ralph 2018). However, one strength of the approach

414    is that the assumption of panmixia within each sub-sample can be relaxed, as our simulated

415    conditions contain stepping stone relationships among demes such that some

416    isolation-by-distance relationships within genetic clusters is accommodated, as long as there is a

417    single source.

418         A second assumption underlying our metric is that a single spatial expansion had actually

419    occurred within an identifiable timeframe. As demonstrated by our simulations with the reduced

420    accuracy as the time since all demes have been colonized ($\tau$) increases, eventually reaching the

421    stationarity conditions of isolation-by-distance, the historical signal of expansion will be erased

422    over long intervals. It is therefore recommended that spatial expansion be tested before

423    attempting to infer the geographic dynamics of an expansion (Excoffier 2004; Wegmann *et al.*

424    2006; Elleouet & Aitken 2018).

425         Key advantages of the GSSA are that no outgroup is required, and that it can be deployed

426    on sparse datasets that are often collected from a variety of non-model species. Notably,

427    correlation to the source of expansion was in general robust to different across-population

428    migration rates ($m$) and deme growth rates—as conditioned by the time allowed for demes to

429    reach their carrying capacity ($\tau_r$). However, accuracy did somewhat decrease with greater

430    migration levels, as one would expect. Also expectedly, the accuracy in identifying the source

431    location was greatly diminished by larger founder sizes ($f$), and the amount of time passed after

432    colonization ended ($\tau$). In the former case, larger $f$ will deflate the founder effect, whereas in the

433    latter case the signal of the founder effect will become erased over time (Nei *et al.* 1975).

434         Regardless of parameter setting, the GSSA approach consistently outperformed the

435    directionality index of Peter and Slatkin's (2013) in all simulations. While such a finding is not

436    surprising given that we violated the latter method's assumption of allelic polarity and used

437    smaller sample sizes than recommended, it highlights the gap filled by our new statistic, which

438    accommodates data typical of non-model organisms. As demonstrated by Peter and Slatkin's

439    (2013) empirical example, in which they found support for the out-of-Africa hypothesis (Li *et al.*

440    2008; DeGiorgio *et al.* 2009) using a comprehensive empirical dataset consisting of over half a

441    million SNPs in more than 1,500 individuals distributed in 55 human populations (Fumagalli *et*

442    *al.* 2011), their statistic is a powerful tool for model organisms. Ours on the other hand, may be

443    best suited for smaller datasets. Furthermore, it is important to note that the two metrics are not

444    precisely inferring the same entity even if the overall objective is to uncover the spatial dynamics

445    of an expansion. Specifically, the smallest RI of the GSSAs is assumed to be associated with the

446    sampling location closest to the actual origin (i.e., first colonized location among those sampled).

447    In contrast, Peter and Slatkin's ψ aims to infer the actual geographical point of expansion, which

448    can be located beyond the sampling localities. This methodological difference reflects our goal

449    of limiting the amount of statistical uncertainty at the expense of decreasing potential inferential

450    capabilities (i.e., limiting the noise introduced by spatial interpolation). Nevertheless, part of our

451    simulation study used sampling locations that included the source, thereby making this

452    comparison between both approaches informative.

453        The utility of our new statistic for non-model organisms is also apparent in the empirical

454    analysis of *Lycaeides* (Fig. 7). These taxa had 8097 and 9074 independent SNPs (i.e., 1 SNP per

455    tag for Alpine *Lycaeides* and  Jackson *Lycaeides* respectively), from which we randomly selected

456    a set of 5000 SNPs for consistency with our simulations. Our approach was able to point, as

457    sources, geographical regions of relatively high predicted suitability according to the ENMs. For

458    Alpine *Lycaeides,* this included a location at the southern end of the range. In contrast, for

459    *Jackson Lycaeides*, the lowest RI suggested an area of marginal LGM suitability (albeit close to

460   regions predicted to be suitable; Fig. 7). That being said, because the inferred entity refers to the

461   sampling location that is closest to the true origin (i.e., first colonized location among those

462   sampled), a perfect match between the ENMs' prediction and the GSSA inference is not

463   expected. Together with the simulation data, these analyses demonstrate that under circumstances

464   of more constrained information and limited sampling relative to those common to model

465   organisms, our approach may be a useful complement to existing methods to infer the geographic

466   history underlying range expansions (Ramachandran *et al.* 2005; François *et al.* 2010; Peter &

467   Slatkin 2013). Our results also highlight how independent lines of evidence can be combined for

468   more robust demographic inference. In contrast to previous attempts to reciprocally validate

469   ENMs and genetic data without spatially explicit inference (Alvarado-Serrano & Knowles 2014),

470   our method allows one to make more explicit comparisons between the hypothesized locations of

471   refugia derived from ENMs and the inferred range expansion dynamics based on genetics. This

472   capability is of particular interest when considering the uncertainty associated with hindcasting

473   using ENMs, as these tools are capable only of generating hypotheses about potential distribution

474   based on abiotic correlates under the assumption of niche conservatism and analogous

475   climate—which is likely to overpredict the real former distribution of species (Soberón 2007;

476   Peterson *et al.* 2010; Alvarado-Serrano & Knowles 2014). By combining these approaches, the

477   confidence in these estimates increases, giving scientists the ability to refine a set of plausible

478   historical scenarios under consideration.


479   **Potential Pitfalls**

23

480         To help improve model-based inferences, several aspects of our approach should be taken

481      into consideration. Besides accurate georeferencing of all samples, the number and spatial

482      distribution of samples should be carefully considered. As it is true for previous methods

483      (Ramachandran *et al.* 2005; Peter & Slatkin 2013), inferences based on the GSSA will be most

484      useful with sampling that maximizes geographical space at the cost of multiple samples per

485      location. Samples clustered in a particular area, and representing only a small portion of the

486      distribution of the taxon of interest, would probably carry limited information and lead to biased

487      estimates. Similarly, our results indicate that precaution should be taken with peripheral localities

488      as they may carry a smaller signal due to noise introduced by boundary effects . In this regard,

489      the advantage of allowing only a single individual per locality should help researchers to more

490      efficiently design their sampling scheme without excessive increases in overall costs, assuming

491      that accessing most of the range of a species is not prohibitive. Furthermore, larger sampling of

492      geographic space could potentially better enable the implementation of a spatial interpolation to

493      identify major colonization routes in more detail (Li & Heap 2011).

494         Like any method using genome-wide SNP data, another relevant consideration is the

495      potential for confounded inference if the sampled SNPs are impacted from parts of the genome

496      under strong natural selection (Sokal *et al.* 1989; Lotterhos & Whitlock 2014), even if they are

497      not linked (Allman & Weissman 2018). For this reason, it would be advisable to test for selective

498      neutrality *a priori*, and to remove loci potentially under selection. An additional consideration is

499      that care should be taken to ensure that inference is not confounded by multiple population

500      histories that involve different sources of expansion. Although use of the GSSA does not

501      explicitly require clustering individuals into populations *a pirori* (Frichot & François 2015;

502    Petkova *et al.* 2016), the conduction of exploratory analyses to obtain some information about

503    population structure may help identify cases of admixture and more than one source for

504    expansion.

505        As a stand-alone summary statistic, the GSSA is itself an exploratory tool to be used with

506    other spatial genomic methods such as EEMS (Petkova *et al.* 2016) , MAPS (Al-Asadi *et al.*

507    2018), or SpaceMix (Bradburd *et al.* 2016): it can  investigate one's data and help formulate

508    model-based hypotheses about population history (House & Hahn 2017). Ideally, the GSSA

509    should be incorporated into an inferential model that allows for testing historical hypotheses and

510    for estimating relevant demographic parameters (He *et al.* 2017a).

**Future Prospects**

512        Further evaluations are needed to assess whether our approach can also detect multiple

513    sources of expansion (a bimodal distribution of raggedness indices could potentially suggest two

514    colonization sources), and if it may allow us to estimate the relative timing of expansions (for

515    instance, by using the slope of the spatial clines in the raggedness of the GSSAs; Fig. 3).

516    Additionally, the GSSA might be used more generally to test for spatial expansion in the first

517    place by using null simulations to determine the significance of spatial clines in the raggedness

518    of GSSAs (Fig. 3). Specifically, the GSSA shape quantified by Harpending's (1994) RI could be

519    used to test for expansion in the context of a null distribution expected under an ahistorical

520    equilibrium between distance and genetic differences (i.e. isolation by distance; IBD). Under the

521    latter scenario, the allelic similarity is expected to be inversely correlated with distance (Wright

522    1943), and hence the spatial cline of raggedness indexes should not be significant. On the

25

523    contrary, under a range expansion this slope is expected to be significantly different from zero.

524    Finally, the GSSA could also be a useful way to explore a broader set of models such as a

525    source/sink population scenarios (Martinez-Solano & Gonzalez 2008), cyclical histories of

526    admixture (Frantz *et al.* 2013; Alvarado-Serrano & Hickerson 2015), or heterogeneous

527    population densities (Excoffier *et al.* 2008). However, the potential for these future directions

528    remains speculative at this point and further work is needed to assess these possibilities.

529          Our approach may also be expanded if there is interest in more precisely identifying the

530    geographic origin of an expansion when the source area is suspected to be missing from the

531    sample. That would require interpolating the expected raggedness index at un-sampled localities

532    to identify the region with the smallest indexes. Implementing an interpolation would also allow

533    us to potentially identify spatially isolated areas of contiguous high raggedness indices, which

534    could serve to uncover instances of expansion from more than one source using an Euclidean

535    allocation algorithm. However, as interpolation is strongly impacted by the number and

536    distribution of observed samples (Stein 2012), its implementation may be not possible when a

537    limited number of locations are available since interpolation accuracy drastically decreases for

538    small sample sizes. Alternatively, a time-difference of arrival approach (TDOA; (Gustafsson *et*

539    *al.* 1994)) could be implemented, as done by Peter and Slatkin's (2013) $\psi$. Commonly used for

540    localization, the TDOA is a triangulation approach that takes advantage of the change in the

541    strength of a signal as distance from the source increases (Gustafsson *et al.* 1994; Drake &

542    Dogancay 2004). Specifically in our case, the pairwise difference in raggedness index between

543    location pairs may be used as the signal change for triangulation to identify the geographic

544    coordinates of the hypothesized source (see eq. 5 in Peter and Slatkin 2013).

545        For those interested in community ecology and species interactions, the GSSA may be

546      used in aggregated geo-referenced population genomic datasets of co-distributed taxa to test for

547      concerted spatio-temporal dynamics between two interacting species (Perkins & Swayne 2001;

548      Wicker *et al.* 2012), the sources of multiple invading species (Sax *et al.* 2007; Johnson *et al.*

549      2009), and the geographic origins of multiple historic domestications (Kanginakudru *et al.* 2008;

550      He *et al.* 2011). Likewise, it may be used to help identify shared regions of secondary contact

551      and hybridization between long-isolated co-distributed pairs of taxa (Remington 1968; Moritz *et*

552      *al.* 2009), or to understand the assembly of whole communities across geographic barriers or

553      trajectories of expansion after global shifts in climate (Avise *et al.* 1987; Ibrahim *et al.* 1996;

554      Avise 2000; Hewitt 2000; Burbrink *et al.* 2016). Additionally, if one were interested in

555      incorporating a resistance surface to correct for effective distances that emerge from landscape

556      features (Spear *et al.* 2010), the GSSA may easily accommodate alternative distance matrices to

557      allow the integration of realistic landscape scenarios.


558    **Conclusion**

559        The GSSA statistic we present here builds on emerging efforts to make phylogeography

560      and population genetics more spatially explicit (Alvarado-Serrano & Hickerson 2015; House &

561      Hahn 2017; Ashander *et al.* 2018; Bradburd *et al.* 2018). By doing so, it improves our ability to

562      estimate the geographical source region (or closest areas) as well as the general direction of

563      range expansions. The use of single samples per location, and un-polarized SNPs, makes this

564      metric well-suited for a wide range of non-model organisms. The GSSA offers not only the

565      capability to serve as an additional spatially explicit summary of genetic patterns for

566    model-based inference when likelihood-based methods are intractable (Beaumont 2010; Pudlo *et*

567    *al.* 2015), but, most importantly, presents a tool for guiding the development of spatially-explicit

568    demographic hypotheses by better accounting for the spatial component of species histories.

569    Further developments of this statistic, including the refinement of an interpolation procedure for

570    geographically sparse samples, should lead to easier incorporation of this tool into

571    simulation-based inferential approaches (Currat *et al.* 2004; Leblois *et al.* 2009; He *et al.* 2017a).

572    We expect that the incorporation of expansion surfaces will enable more complex scenarios that

573    include environmental heterogeneity and explicit geographic barriers to be modeled into these

574    simulation approaches (Ray & Excoffier 2010; Joseph *et al.* 2016). Likewise, we expect clines in

575    GSSA estimates to be able to help identify relevant spatial features, such as shared contact zones

576    between populations that have been colonized from different expansion clusters (Swenson 2010),

577    or consistent geographic barriers that have maintained isolation of populations during expansions

578    and hence have promoted differentiation (Carstens *et al.* 2005; Potter *et al.* 2017). The GSSA

579    offers a useful and flexible tool that complements existing methods for improved understanding

580    of the processes governing population differentiation and spatial patterns of genomic diversity.

581    ## Software

582    A program to calculate the GSSA and Harpending's (1944) raggedness index from

583    geo-referenced genotypes and associated user-defined geographic distances among sampled

584    locations is available in https://bitbucket.org/diegofalvarado-s/GSSA_v0.0. This software also

585    contains the pipeline needed to recreate the simulations we used to test the accuracy of the GSSA

586 and Peter and Slatkin's $\psi$. A full step-by-step implementation of our method, including

587 comments and visualization of the intermediate outputs, is available at the following jupyter

588 notebook: https://mybinder.org/v2/gh/ftempo/GSSA/master.

## Acknowledgements

# References

Al-Asadi, H., Petkova, D., Stephens, M., & Novembre, J. (2018). Estimating recent migration and population size surfaces. *bioRxiv*, **doi:**10.1101/365536.

Allman, B.E., & Weissman, D.B. (2018). Hitchhiking in space: Ancestry in adapting, spatially extended populations. *Evolution*, **72**, 722–734.

Alvarado-Serrano, D.F., & Hickerson, M.J. (2015). Spatially explicit summary statistics for historical population genetic inference. *Methods in Ecology and Evolution*, **7**, 418–427.

Alvarado-Serrano, D.F., & Knowles., L.L. (2014). Ecological niche models in phylogeographic studies: Applications, advances and precautions. *Molecular Ecology Resources*, **14**, 233–248.

Ashander, J., Ralph, P., McCartney-Melstad, E., & Shaffer, H.B. (2018), Demographic inference in a spatially-explicit ecological model from genomic data: A proof of concept for the Mojave Desert Tortoise. *bioRxiv*. **doi:** 10.1101/354530.

Avise, J.C. (2000). *Phylogeography: The history and formation of species*. Harvard University Press. Cambridge.

Avise, J.C., Arnold, J., Ball, R.M., Bermingham, E., Lamb, T., Neigel, J.A., …, Saunders, N.C. (1987). Intraspecific phylogeography: The mitochondrial DNA bridge between population genetics and systematics. *Annual Review of Ecology and Systematics*, **18**, 489–522.

Beaumont, M. (2010). Approximate Bayesian computation in evolution and ecology. *Annual Review of Ecology, Evolution, and Systematics*, **41**, 379–406.

Braconnot, P., Otto-Bliesner, B., Harrison, S., Joussaume, S., Peterchmitt, J.-Y., Abe-Ouchi, A., …, Zhao, Y. (2007). Results of PMIP2 coupled simulations of the Mid-Holocene and Last Glacial Maximum--Part 1: experiments and large-scale features. *Climate of the Past*, **3**, 261–277.

Bradburd, G.S., Coop, G.M., & Ralph, P.L. (2018). Inferring continuous and discrete population genetic structure across space. *Genetics*, **210**, 33–52.

Bradburd, G.S., Ralph, P.L., & Coop, G.M. (2016). A spatial framework for understanding population structure and admixture. *PLoS genetics*, **12**, e1005703.

Branco, C., Velasco, M., Benguigui, M., Currat, M., Ray, N., & Arenas, M. (2018). Consequences of diverse evolutionary processes on american genetic gradients of modern humans. *Heredity,* **doi:** 10.1038/s41437-018-0122-x.

Brown,J.L., Weber, J.J., Alvarado-Serrano, D.F., Hickerson, M.J., Franks, S.J., & Carnaval, A.C.

632 (2016). Predicting the genetic consequences of future climate change: The power of
633 coupling spatial demography, the coalescent, and historical landscape changes. *American*
634 *Journal of Botany*, **103**, 153–163.

635 Burbrink, F.T., Chan, Y.L., Myers, E.A., Ruane, S., Smith, B.T., & Hickerson, M.J. (2016).
636 Asynchronous demographic responses to Pleistocene climate change in Eastern Nearctic
637 vertebrates. *Ecology Letters*, **19**, 1457–1467.

638 Cardozo, A.L.P., Farias, E.G.G., Rodrigues-Filho, J.L., Moteiro, I.B., Scandolo, T.M. & Dantas,
639 T.V. (2018). Feeding ecology and ingestion of plastic fragments by *Priacanthus arenatus*:
640 What's the fisheries contribution to the problem? *Marine Pollution Bulletin*, **130**, 19–27.

641 Carstens, B.C., Brunsfeld, S.J., Demboski, J.R., Good, J.D., & Sullivan, J. (2005). Investigating
642 the evolutionary history of the Pacific Northwest mesic forest ecosystem: Hypothesis
643 testing within a comparative phylogeographic framework. *Evolution*, **59**, 1639–1652.

644 Charles H., & Dukes, J.S. (2008). Impacts of invasive species on ecosystem services. In:
645 *Biological Invasions Ecological Studies*, pp. 217–237. Springer, Heidelberg.

646 Cristescu, M.E. (2015). Genetic reconstructions of invasion history. *Molecular Ecology*, **24**,
647 2212–2225.

648 Currat, M., Ray, N., & Excoffier, L. (2004). Splatche: a program to simulate genetic diversity
649 taking into account environmental heterogeneity. *Molecular Ecology Notes*, **4**, 139–142.

650 Davis, C.D., Epps, C.W., Flitcroft, R.L., & Bank,s M.A. (2018). Refining and defining riverscape
651 genetics: How rivers influence population genetic structure. *Wiley Interdisciplinary*
652 *Reviews: Water*, **5**, e1269.

653 DeGiorgio, M., Jakobsson, M., & Rosenberg, N.A. (2009). Explaining worldwide patterns of
654 human genetic variation using a coalescent-based serial founder model of migration
655 outward from Africa. *Proceedings of the National Academy of Sciences*, **106**, 16057–16062.

656 Drake, S.R., & Dogancay, K. (2004). Geolocation by time difference of arrival using hyperbolic
657 asymptotes. In: *2004 IEEE International Conference on Acoustics, Speech, and Signal*
658 *Processing*, pp. 361–364.

659 Elleouet, J.S., & Aitken, S.N. (2018). Exploring Approximate Bayesian Computation for
660 inferring recent demographic history with genomic markers in non-model species.
661 *Molecular Ecology Resources*, **doi:** 10.1111/1755-0998.12758.

662 Excoffier, L. (2004). Patterns of DNA sequence diversity and genetic structure after a range
663 expansion: Lessons from the infinite-island model. *Molecular Ecology*, **13**, 853–864.

664 Excoffier, L., Dupanloup, I., Huerta-Sánchez, E., Sousa, V.C., & Foll, M. (2013). Robust

31

demographic inference from genomic and SNP data. *PLoS genetics*, **9**, e1003905.

Excoffier, L., & Foll, M. (2011). Fastsimcoal: A continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics*, **27**, 1332–1334.

Excoffier, L., Foll, M., & Petit, R.J. (2008). Genetic consequences of range expansions. *Annual Review of Ecology, Evolution, and Systematics*, **40**, 481–501.

Excoffier, L., Hofer, T., & Foll, M. (2009). Detecting loci under selection in a hierarchically structured population. *Heredity*, **103**, 285–298.

Fagua, J.C., & Gonzalez, V.H. (2007). Growth rates, reproductive phenology, and pollination ecology of *Espeletia grandiflora* (Asteraceae), a giant Andean caulescent rosette. *Plant Biology* , **9**, 127–135.

Fordham, D.A., Brook, B.W., Moritz, C., & Nogués-Bravo, D. (2014). Better forecasts of range dynamics using genetic data. *Trends in Ecology & Evolution*, **29**, 436–443.

Forister, M.L., Gompert, Z., Fordyce, J.A., & Nice, C.C. (2011). After 60 years, an answer to the question: What is the Karner blue butterfly? *Biology Letters*, **7**, 399–402.

Fraimout, A., Debat, V., Fellous, S., Hufbauer, R.A., Foucaud, J., Pudlo, P., …, Estoup, A. (2017). Deciphering the routes of invasion of *Drosophila suzukii* by means of ABC Random Forest. *Molecular Biology and Evolution*, **34**, 980–996.

François., O., Currat, M, Ray, N., Han, E., Excoffier, L. & Novembre, J. (2010). Principal component analysis under population genetic models of range expansion and admixture. *Molecular Biology and Evolution*, **27**, 1257–1268.

Frantz, L.A.F., Schraiber, J.G., Madsen, O., Megens, H.J., Bosse, M., Paudel, Y., …, Groenen, M.A. (2013). Genome sequencing reveals fine scale diversification and reticulation history during speciation in *Sus*. *Genome Biology*, **14**, R107.

Frichot, E., & François, O. (2015). LEA: An R package for landscape and ecological association studies. *Methods in Ecology and Evolution*, **6**, 925–929.

Fumagalli, M., Sironi, M., Pozzoli, U., Ferrer-Admetlla, A., Pattini, L. & Nielsen, R. (2011). Signatures of environmental genetic adaptation pinpoint pathogens as the main selective pressure through human evolution. *PLoS genetics*, **7**, e1002355.

Galante, P.J., Alade, B., Muscarella, R., Jansa, S.A., Goodman, S.M., & Anderson, R.P. (2018). The challenge of modeling niches and distributions for data-poor species: a comprehensive approach to model complexity. *Ecography*, **41**, 726–736.

Gaston KJ (2003) *The Structure and Dynamics of Geographic Ranges*. Oxford University Press. Oxford.

Gompert, Z., Comeault, A.A., Farkas, T.E., Feder, J.L., Parchman, T.L., Buerkle, C.A., & Nosil, P. (2014a). Experimental evidence for ecological selection on genome variation in the wild. *Ecology Letters*, **17**, 369–379.

Gompert, Z., Fordyce, J.A., Forister, M.L., Shapiro, A.M., & Nice, C.C. (2006). Homoploid hybrid speciation in an extreme habitat. *Science*, **314**, 1923–1925.

Gompert, Z., Forister, M.I., Fordyce, J.A., Nice, C.C., Williamson, R.J., & Buerkle, C.A. (2010). Bayesian analysis of molecular variance in pyrosequences quantifies population genetic structure across the genome of *Lycaeides* butterflies. *Molecular Ecology*, **19**, 2455–2473.

Gompert, Z., Lucas, L.K., Buerkle, C.A., Forister, M.L., Fordyce, J.A., & Nice, C.C. (2014b). Admixture and the organization of genetic diversity in a butterfly species complex revealed through common and rare genetic variants. *Molecular Ecology*, **23**, 4555–4573.

Gustafsson, F., Gunnarsson, S., & Ljung, L. (1994). On time-frequency resolution of signal properties using parametric techniques. In: 1994 *Proceedings of 33rd IEEE Conference on Decision and Control*, pp. 2259–2264.

Gutenkunst, R.N., Hernandez, R.D., Williamson, S.H., Bustamante, C.D. (2009). Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS genetics*, **5**, e1000695.

Hallatschek, O., Hersen, P., Ramanathan, S., & Nelson, D.R. (2007) Genetic drift at expanding frontiers promotes gene segregation. *Proceedings of the National Academy of Sciences of the United States of America*, **104**, 19926–19930.

Hallatschek, O., & Nelson DR (2008) Gene surfing in expanding populations. *Theoretical Population Biology*, **73**, 158–170.

Harpending, H.C. (1994). Signature of ancient population growth in a low-resolution mitochondrial DNA mismatch distribution. *Human Biology*, **66**, 591–600.

Henn, B.M., Cavalli-Sforza, L.L., & Feldman, M.W. (2012). The great human expansion. *Proceedings of the National Academy of Sciences of the United States of America*, **109**, 17758–17764.

He, Q., Prado, J.R., & Knowles, L.L. (2017). Inferring the geographic origin of a range expansion: Latitudinal and longitudinal coordinates inferred from genomic data in an ABC framework with the program x-origin. *Molecular Ecology*, **26**, 6908–6920.

Hewitt, G. (2000). The genetic legacy of the Quaternary ice ages. *Nature*, **405**, 907–913.

33

730  He, Z., Zhai, W., Wen, H., Tang, T., Wang, Y., Lu, X., …, Shi,S. (2011). Two evolutionary
731      histories in the genome of rice: The roles of domestication genes. *PLoS genetics*, **7**,
732  e1002100.

733  Hijmans, R.J. (2012). Cross-validation of species distribution models: Removing spatial sorting
734      bias and calibration with a null model. *Ecology*, **93**, 679–688.

735  House, G.L., & Hahn, M.W. (2018). Evaluating methods to visualize patterns of genetic
736      differentiation on a landscape. *Molecular Ecology Resources*, **18**, 448-460.

737  Ibrahim, K.M., Nichols, R.A., & Hewitt, G.M. (1996). Spatial patterns of genetic variation
738      generated by different forms of dispersal during range expansion. *Heredity*, **77**, 282–291.

739  Johnson, P.T.J., Olden, J.D., Solomon, C.T., Vander Zanden, M.J. (2009). Interactions among
740      invaders: Community and ecosystem effects of multiple invasive species in an experimental
741  aquatic system. *Oecologia*, **159**, 161–170.

742  Jones, E., Oliphant, T., & Peterson, P. (2001). SciPy: Open source scientific tools for Python,
743      **doi:** 10.1234/12345678.

744  Joseph, T.A., Hickerson, M.J., & Alvarado-Serrano, D.F. (2016). Demographic inference under a
745      spatially continuous coalescent model. *Heredity*, **117**, 94–99.

746  Kanginakudru, S., Metta, M., Jakati, R.D., & Nagaraju, J. (2008). Genetic evidence from Indian
747      red jungle fowl corroborates multiple domestication of modern day chicken. *BMC
748  Evolutionary Biology*, **8**, 174.

749  Kimura, M., & Weiss, G.H. (1964). The stepping stone model of population structure and the
750      decrease of genetic correlation with distance. *Genetics*, **49**, 561–576.

751  Kirkpatrick M., & Barton, N.H. (1997). Evolution of a species' range. *The American naturalist*,
752      **150**, 1–23.

753  Knowles, L.L., & Alvarado-Serrano, D.F. (2010). Exploring the population genetic consequences
754      of the colonization process with spatio-temporally explicit models: Insights from coupled
755  ecological, demographic and genetic models in montane grasshoppers. *Molecular Ecology*,
     **19**, 3727–3745.

757  Konečný, A., Estoup, A., Duplantier, J.-M., Bryja, J., Ba, K., Galan, M., …, Cosson, J.F. (2013).
758      Invasion genetics of the introduced black rat (*Rattus rattus*) in Senegal, West Africa.
759  *Molecular Ecology*, **22**, 286–300.

760  Leblois, R., Estoup, A., Rousset, F. (2009). IBDSim: A computer program to simulate genotypic
761      data under isolation by distance. *Molecular Ecology Resources*, **9**, 107–109.

762    Legendre P., & Legendre, L.F.J. (2012). *Numerical Ecology*. Elsevier. Philadelphia.

763    Li, J.Z., Absher, D.M., Tang, H., Southwick, A.M., Casto, A.M., Ramachandran, S., …, Myers,
764        R.M. (2008). Worldwide human relationships inferred from genome-wide patterns of
765        variation. *Science*, **319**, 1100–1104.

766    Li J., & Heap, A.D. (2011). A review of comparative studies of spatial interpolation methods in
767        environmental sciences: Performance and impact factors. *Ecological Informatics*, **6**,
768        228–241.

769    Lotterhos, K.E., & Whitlock, M.C. (2014). Evaluation of demographic history and neutral
770        parameterization on the performance of $F_{ST}$ outlier tests. *Molecular Ecology*, **23**,
771        2178–2192.

772    Lundgren, E., & Ralph, P.L. (2018). Are populations like a circuit? The relationship between
773        isolation by distance and isolation by resistance. *bioRxiv*, **doi:** 10.1101/451328.

774    Martinez-Solano, I., Gonzalez, E.G. (2008). Patterns of gene flow and source-sink dynamics in
775        high altitude populations of the common toad *Bufo bufo* (Anura: Bufonidae). *Biological
776        Journal of the Linnean Society*, **95**, 824–839.

777    Menozzi, P., Piazza, A., & Cavalli-Sforza, L. (1978). Synthetic maps of human gene frequencies
778        in Europeans. *Science*, **201**, 786–792.

779    Miraldo, A., Hewitt, G.M., Paulo, O.S., & Emerson, B.C. (2011). Phylogeography and
780        demographic history of *Lacerta lepida* in the Iberian Peninsula: Multiple refugia, range
781        expansions and secondary contact zones. *BMC Evolutionary Biology*, **11**, 170.

782    Moritz, C., Hoskin, C.J., MacKenzie, J.B., Phillips, B.L., Tonione, M., Silva, N., …, Graham,
783        C.H. (2009). Identification and dynamics of a cryptic suture zone in tropical rainforest.
784        *Proceedings of the Royal Society B: Biological Sciences*, **276**, 1235–1244.

785    Muscarella, R., Galante, P.J., Soley-Guardia, M., Boria, R.A., Kass, J.M., Uriarte, M., &
786        Anderson, R.P. (2014). ENMeval: An R package for conducting spatially independent
787        evaluations and estimating optimal model complexity for Maxent ecological niche models.
788        *Methods in Ecology and Evolution*, **5**, 1198–1205.

789    Nei, M., Maruyama, T., & Chakraborty, R. (1975). The bottleneck effect and genetic variability
790        in populations. *Evolution*, **29**, 1–10.

791    Nice, C.C., Gompert, Z., Fordyce, J.A., Forister, M.L., Lucas, L.K., & Buerkle, C.A. (2013).
792        Hybrid speciation and independent evolution in lineages of alpine butterflies. *Evolution*, **67**,
793        1055–1068.

794    Novembre, J., & Stephens, M. (2008). Interpreting principal component analyses of spatial
795        population genetic variation. *Nature Genetics*, **40**, 646–649.

796    Patterson, N., Price, A.L., & Reich, D. (2006). Population structure and eigenanalysis. *PLoS
797        Genetics*, **2**, e190.

798    Peischl, S., Dupanloup, I., Kirkpatrick, M., & Excoffier, L. (2013). On the accumulation of
799        deleterious mutations during range expansions. *Molecular Ecology*, **22**, 5972–5982.

800    Pemberton, T.J., DeGiorgio, M., & Rosenberg, N.A. (2013). Population structure in a
801        comprehensive genomic data set on human microsatellite variation. *G3: Genes, Genomes,
802        Genetics*, **3**, 891–907.

803    Perkins, L.E., & Swayne, D.E. (2001). Pathobiology of A/chicken/Hong Kong/220/97 (H5N1)
804        avian influenza virus in seven gallinaceous species. *Veterinary Pathology*, **38**, 149–164.

805    Peter, B.M., & Slatkin, M. (2013). Detecting range expansions from genetic data. *Evolution*, **67**,
806        3274–3289.

807    Peter, B.M., & Slatkin, M. (2015). The effective founder effect in a spatially expanding
808        population. *Evolution*, **69**, 721–734.

809    Peterson, K.R., Pfister, D.H., & Bell, C.D. (2010). Cophylogeny and biogeography of the fungal
810        parasite *Cyttaria* and its host *Nothofagus*, southern beech. *Mycologia*, **102**, 1417–1425.

811    Petkova, D., Novembre, J., & Stephens, M. (2016). Visualizing spatial population structure with
812        estimated effective migration surfaces. *Nature Genetics*, **48**, 94–100.

813    Pharo, E.J., & Zartman, C.E. (2007). Bryophytes in a changing landscape: The hierarchical
814        effects of habitat fragmentation on ecological and evolutionary processes. *Biological
815        Conservation*, **135**, 315–325.

816    Phillips, S.J., Anderson, R.P., & Schapire, R.E., (2006). Maximum entropy modeling of species
817        geographic distributions. *Ecological Modelling*, **190**, 231–259.

818    Phillips, B.L., Kelehear, C., Pizzatto, L., Brown, G.P., Barton, D., & Shine, R. (2010). Parasites
819        and pathogens lag behind their host during periods of host range advance. *Ecology*, **91**,
820        872–881.

821    Pierce, J.L., Meyer, G.A., & Jull, A.J.T. (2004). Fire-induced erosion and millennial-scale
822        climate change in northern ponderosa pine forests. *Nature*, **432**, 87–90.

823    Pires, T.H.S., Farago, T.B., Campos, D.F., Cardoso, G.M., & Zuanon, J. (2016). Traits of a
824        lineage with extraordinary geographical range: Ecology, behavior and life-history of the
825        sailfin tetra *Crenuchus spilurus*. *Environmental Biology of Fishes*, **99**, 925–937.

826 Potter, S., Xue, A.T., Bragg, J.G., Rosauer, D.F., Roycroft, E.J., & Moritz, C. (2017).
827 Pleistocene climatic changes drive diversification across a tropical savanna. *Molecular*
828 *Ecology*, **27**, 520–532.

829 Pudlo, P., Marin, J.-M., Estoup, A., Cornuet, J.-M., Gautier, & M., Robert, C.P. (2015). Reliable
830 ABC model choice via random forests. *Bioinformatics* **32**, 859-866.

831 Radosavljevic, A., & Anderson, R.P. (2014) Making better Maxent models of species
832 distributions: complexity, overfitting and evaluation. *Journal of Biogeography*, **41**,
833 629–643.

834 Ramachandran, S., Deshpande, O., Roseman, C.C., Rosenberg, N.A., Feldman, M.W.,
835 Cavalli-Sforza, L. (2005). Support from the relationship of genetic and geographic distance
836 in human populations for a serial founder effect originating in Africa. *Proceedings of the*
837 *National Academy of Sciences of the United States of America*, **102**, 15942–15947.

838 Ramírez-García, P., López-Blanco J., & Ocaña, D. (1998). Mangrove vegetation assessment in
839 the Santiago River Mouth, Mexico, by means of supervised classification using LandsatTM
840 imagery. *Forest Ecology and Management*, **105**, 217–229.

841 Ray, N., & Excoffier, L. (2010). A first step towards inferring levels of long-distance dispersal
842 during past expansions. *Molecular Ecology Resources*, **10**, 902–914.

843 Remington, C.L. (1968). Suture-zones of hybrid interaction between recently joined biotas. In:
844 *Evolutionary Biology*, pp. 321–428. Springer, Boston.

845 Roberts, D.R., & Hamann, A. (2015). Glacial refugia and modern genetic diversity of 22 western
846 North American tree species. *Proceedings of the Royal Society B: Biological Sciences*, **282**,
847 20142903.

848 Ruiz-Cooley, R.I., Ballance, L.T., & McCarthy, M.D. (2013). Range expansion of the jumbo
849 squid in the NE Pacific: δ15N decrypts multiple origins, migration and habitat use. *PloS*
850 *one*, **8**, e59651.

851 Sax, D.F., Stachowicz, J.J., Brown, J.H., Bruno, J.F., Dawson, M.N, Gaines, S.D., …, Rice, W.R.
852 (2007). Ecological and evolutionary insights from species invasions. *Trends in Ecology &*
853 *Evolution*, **22**, 465–471.

854 Schrider, D.R., & Kern, A.D. (2018). Supervised machine learning for population genetics: A
855 new paradigm. *Trends in Genetics,* **34**, 301-312.

856 Scott, D.W. (1979). On optimal and data-based histograms. *Biometrika*, **66**, 605–610.

857 Shirk, A.J., & Cushman, S.A. (2014). Spatially-explicit estimation of Wright's neighborhood size

in continuous populations. *Frontiers in Ecology and Evolution*, **2**, 177.

Slatkin, M. (1987) Gene flow and the geographic structure of natural populations. *Science*, **236**, 787–792.

Slatkin, M., & Excoffier, L. (2012). Serial founder effects during range expansion: A spatial analog of genetic drift. *Genetics*, **191**, 171–181.

Smouse, P.E., & Peakall, R. (1999). Spatial autocorrelation analysis of individual multiallele and multilocus genetic structure. *Heredity*, **82**, 561–573.

Soberón, J. (2007) Grinnellian and Eltonian niches and geographic distributions of species. *Ecology Letters*, **10**, 1115–1123.

Sokal, R.R., Harding, R.M., & Oden, N.L. (1989). Spatial patterns of human gene frequencies in Europe. *American Journal of Physical Anthropology*, **80**, 267–294.

Spear, S.F., Balkenhol, N., Fortin, M.-J., McRae, B.H., & Scribner, K. (2010). Use of resistance surfaces for landscape genetic studies: considerations for parameterization and analysis. *Molecular Ecology*, **19**, 3576–3591.

Stein, M.L. (2012). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer. New York.

Sturges, H.A. (1926). The choice of a class interval. *Journal of the American Statistical Association*, **21**, 65–66.

Swenson, N.G. (2010). Mapping the suturing of a continental biota. *Molecular Ecology*, **19**, 5324–5327.

Thompson, R.S., & Anderson, K.H. (2000). Biomes of western North America at 18,000, 6000 and 0 14C yr bp reconstructed from pollen and packrat midden data. *Journal of Biogeography*, **27**, 555–584.

Thompson, L.G., Mosley-Thompson, E., Davis, M. Lin, P.N., Yao, T., Dyurgerov, M., & Dai, J. (1993). "Recent warming": ice core evidence from tropical ice cores with emphasis on Central Asia. *Global and Planetary Change*, **7**, 145–156.

Thuiller, W., Albert, C., Araújo, M.B., Berry, P.M., Cabeza, M., Guisan, A., …, Zimmermann, E. (2008). Predicting global change impacts on plant species' distributions: future challenges. *Perspectives in Plant Ecology, Evolution and Systematics*, **9**, 137–152.

Wegmann, D., Currat, M., & Excoffier, L. (2006). Molecular diversity after a range expansion in heterogeneous environments. *Genetics*, **174**, 2009–2020.

Wicker, E., Lefeuvre, P., de Cambiaire, J.-C., Lemaire, C., Poussier, S., & Prior, P. (2012).

889      Contrasting recombination patterns and demographic histories of the plant pathogen

890      *Ralstonia solanacearum* inferred from MLSA. *The ISME Journal*, **6**, 961–974.

891   Wright, S. (1943). Isolation by distance. *Genetics*, **28**, 114–138.

## Data Accessibility

893   The whole pipeline and associated files needed to recreate the simulations used is available at the

894   following bitbucket repository: https://bitbucket.org/diegofalvarado-s/GSSA_v0.0

## Author Contributions

896      DFA and MJH conceived the GSSA statistic and DFA led the development of it. DFA coded

897   the simulations, analyses and produced the figures. DFA and MJH wrote the manuscript.

# Tables

898

**Table 1.** Parameters used in simulations. Parameter values in bold were used for all simulation 900experiments except those in which the parameter is being tested.

899

| Parameter | Description | Values tested | Unit |
|---|---|---|---|
| $\tau$ | Time allowed for the simulation to continue after last deme has been colonized and reached its carrying capacity | **0.0** 0.5 1.0 | Coalescent units |
| $f$ | Proportion of individuals of a source deme that settle each uncolonized deme | **0.002** 0.01 0.1 | Proportion of $N_K$ |
| $\tau_r$ | Time allowed for demes to reach $N_K$ after being colonized | **0.01** 0.05 0.10 | Coalescent units |
| $m$ | Pairwise probability of migration per individual per generation between neighboring colonized demes | **0.000** 0.001 0.010 | Proportion |
| $N_K$ | Effective population size of deme at carrying capacity | **1000** | Number of individual |
| $\tau_c$ | Time lag between source deme reaching $N_K$ and colonizing a new deme | **0.01** | Coalescent units |

40

## Figure Legends

**Fig. 1.** Schematic cartoon describing the construction of the geographic spectrum of shared alleles (GSSA) across five diploid genotypes that are each sampled from a unique location. Location-specific genotype matrices ($G_i$), which capture the presence/absence of minor alleles at each location, locus and DNA strand (**panel a**), and cartoon representation of location-specific vectors ($\vec{d_i}$), which capture the geographic distances between each focal location and all of the other locations in the sample (**panel b**). Based on the $G_i$ matrices and $\vec{d_i}$ distance vectors, corresponding vectors ($\vec{S_i}$) are constructed from the aggregated relative spatial distribution of minor alleles for each location (**panel c**). To facilitate interpretation, the color of each $\vec{S_i}$ element follows the location colors in panel a, and indicates their derivation with respect to each one of the sampled locations. The $\vec{S_i}$ vectors are then converted into corresponding geo-genetic histograms ($h_{gen_i}$), by applying a common binning scheme based on Sturge's (1926) equation (**panel d**). For each location, a corresponding geographic histogram ($h_{geo_i}$) is constructed, using the same binning scheme as before, from the geographic-distances separating the location from all other sampled locations (**panel e**). The number of observations for corresponding distance classes in both histograms, $h_{gen_i}$ and $h_{geo_i}$, are then regressed against each other for each location (**panel f**). From these regressions, the GSSA vector for each individual is built by taking the absolute value of the regression residuals for each histogram's distance class (**panel g**).

**Fig. 2.** Schematic of the sequential range expansion models used to quantify the spatial association of the GSSA and expansion source location. Each panel correspond to simulations started from a different source (marked in red and denoted by letter). Demes are shaded according to their colonization time, with darker colors corresponding to earlier colonization times. Sampled demes are surrounded by a square.

**Fig. 3.** Association between time at which each sampled locality is colonized and the raggedness index calculated from the GSSA at each of these localities under all four simulated scenarios. Colonization time step denotes the time at which a locality was colonized (with the origin always being zero). The impact of each parameter on the ability to recover the simulated colonization history is evidenced by changes in the slope of the association. Different colors are used for each parameter value for increased visibility.

**Fig. 4.** Probability across 1000 simulation replicates of the the raggedness index calculated from the GSSA vector for identifying the source deme given the four serial range expansion scenarios and model parameter values. Higher probabilities of correct source identifications are a function of darker shadings. See Table 1 for model parameter definitions.

**Fig. 5.** Geographic position of the suggested sources of expansion according to the raggedness index calculated from the GSSA vector (left column) and directionality index (right column). A red circle denotes the true source used for each simulation scenario, whereas black circles indicate the locations sampled (shown only for the GSI-RI approach). Results presented are

41

938  aggregated across 1000 simulations, with the relative frequency that each position was identified
939 as the source of expansion indicated by the color hue (darker colors corresponding to greater
940 frequencies). Note that in the vast majority of instances, the GSSA-RI approach either correctly
941 infers the source or identifies a close neighbor location.

942  **Fig. 6.** Comparative accuracies of the raggedness index calculated from the GSSA and
943 directionality indexin identifying the source deme for serial range expansions. The mean
944 accuracy (a) per model parameter tested and corresponding root mean square error (RMSE) (b)
945 obtained under both methods (GSSA-RI in blue, directionality index in red) is shown under
946 different model parameter values.

947  **Fig. 7.** Comparison of LGM hindcasts of range distributions obtained with Ecological Niche
948 Models (ENMs). The grey circles indicate the sampling localities on which the inferences are
949 based. The suggested sources of expansion according to the GSSA-RI approach are indicated by
950 a green rhombus with a "G" on it. Warmer colors indicate greater potential suitability as
951 estimated by the ENM. a) Model and inferences for Alpine *Lycaeides*; b) model and inference
952 for Jackson *Lycaeides*. Refugia, identified as the top 10th percentile most suitable areas, are
953 denoted by black polygons.

**Supp. Table 1.** Tuning results for ecological niche models of *Lycaeides* taxa based on jacknife cross-validation. Abbreviations: AUC: area-under the ROC curve; OR.10: 10th percentile omission rate; AICc: corrected Akaike Information Criterion; L: linear features, Q: quadratic features, LQ: linear-quadratic features; H: hinge features. Models for each species are ordered according to their ranking, with the best model on top.

**Supp. Fig. 1.** Estimation accuracy of the raggedness index calculated from the GSSA vector under alternative histogram binning schemes. Accuracy of the GSSA-based inference when constructed using Sturge's binning equation is compared against its accuracy when constructed using a coarser (a), or a finer-grain binning scheme (b). Each dot corresponds to the mean accuracy with which the source of expansion is inferred for each parameter combination and binning scheme used (points are colored according to the parameters used in the simulations; see Table 1). The dashed line corresponds to a 1:1 line.

**Supp. Fig. 2.** Example of GSSA vectors for (a) last-colonized and (b) the expansion source locality in one of the simulations. The frequency of each of the GSSA bin ($b_k$) is depicted along the corresponding raggedness index (Simulation parameters used: $\tau = 0$; $f = 0.002$; $\tau_r = 0.01$; $m = 0$; $N_K = 1000$; $\tau_c = 0.01$; simulated source = "origin B").

**Supp. Fig. 3.** Geographic position of the sampled locality closest to the source of expansion according to the GSSA-RI approach. A red circle denotes the locality, among those sampled (indicated with black incircles), that was first colonized (note the simulated source itself, denoted by a black square, is not sampled). Results presented are aggregated across simulations, with the relative frequency that each position was identified as the first to be colonized indicated by the color hue (darker colors corresponding to greater frequencies). Note that in the vast majority of instances, the GSSA-RI approach correctly infers the closest neighbor to the source or another closeby neighbor location.
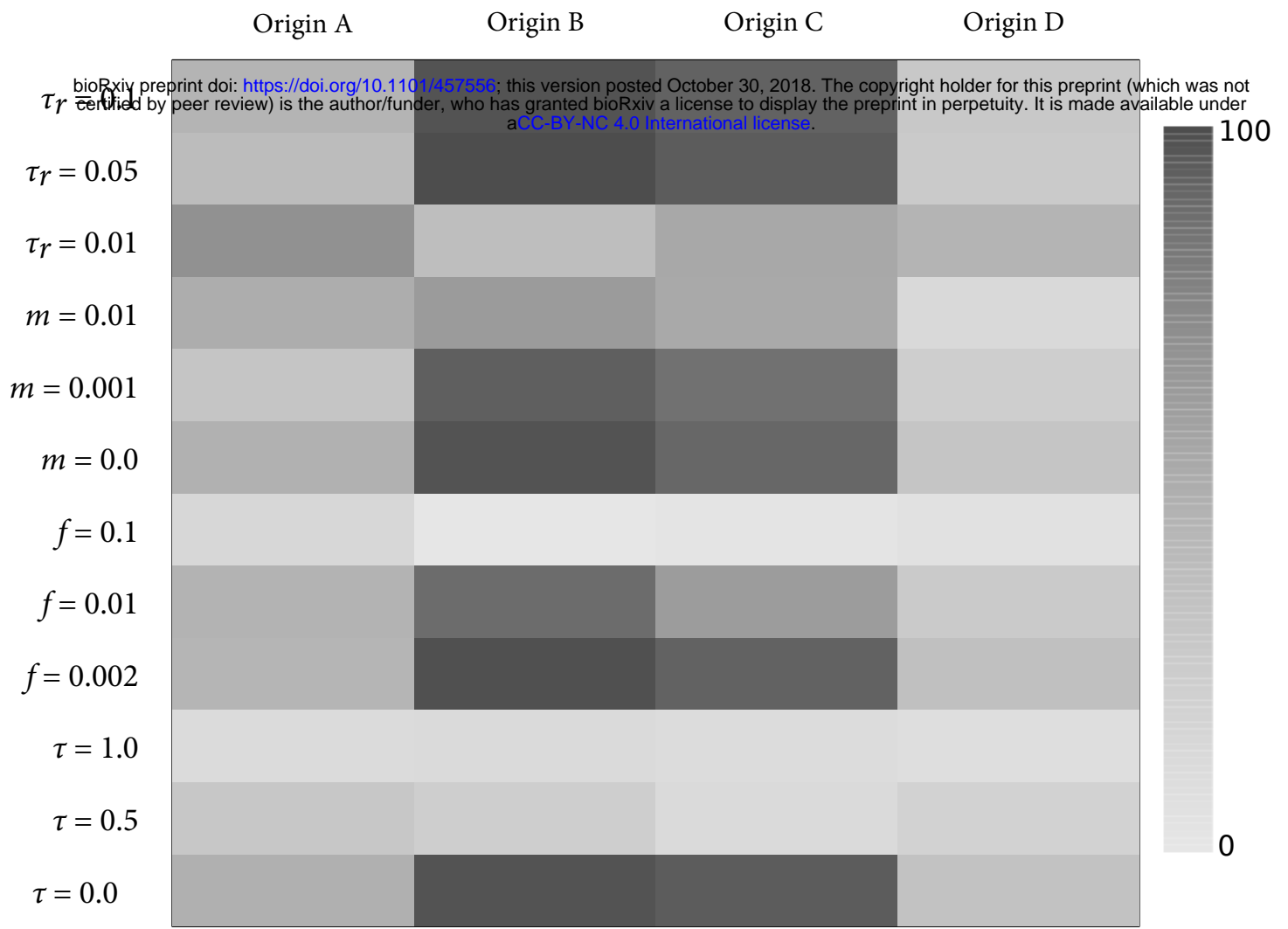
Fig. 1

**Fig. 2**

Fig. 3a

Fig. 3b

**Fig. 4**

PSI



Fig. 5

**Fig. 6**

a)

b)

Fig. 7

| Taxon | Features | Regularization Multiplier | Mean OR.10 | Mean AUC | AICc | No. parameters |
|---|---|---|---|---|---|---|
| Alpine | L | 1 | 0.2 | 0.89182 | 206.6493599 | 3 |
| Alpine | L | 1.5 | 0.2 | 0.89661 | 203.2155179 | 2 |
| Alpine | H | 1 | 0.2 | 0.91757 | - | 10 |
| Alpine | LQ | 2.5 | 0.2 | 0.92362 | 192.4424708 | 2 |
| Alpine | LQ | 3 | 0.2 | 0.92511 | 191.7289706 | 1 |
| Alpine | LQ | 3.5 | 0.2 | 0.92638 | 192.100387 | 1 |
| Alpine | LQ | 4 | 0.2 | 0.926675 | 192.481438 | 1 |
| Alpine | LQ | 2 | 0.2 | 0.942205 | 189.4175236 | 2 |
| Alpine | LQ | 1.5 | 0.2 | 0.953875 | 186.7180064 | 2 |
| Alpine | LQ | 1 | 0.2 | 0.954975 | 188.7346517 | 3 |
| Alpine | L | 4 | 0.3 | 0.894535 | 204.4186749 | 2 |
| Alpine | L | 3.5 | 0.3 | 0.894875 | 204.1054873 | 2 |
| Alpine | L | 3 | 0.3 | 0.896515 | 203.8267215 | 2 |
| Alpine | L | 2 | 0.3 | 0.896815 | 203.3794817 | 2 |
| Alpine | L | 2.5 | 0.3 | 0.89716 | 203.5840395 | 2 |
| Alpine | H | 1.5 | 0.3 | 0.919485 | - | 9 |
| Alpine | LQ | 0.5 | 0.3 | 0.948905 | 191.9917093 | 4 |
| Alpine | H | 3.5 | 0.4 | 0.900645 | 213.0286531 | 4 |
| Alpine | H | 3 | 0.4 | 0.904335 | 210.8766708 | 4 |
| Alpine | H | 2 | 0.4 | 0.921855 | 254.9666543 | 7 |
| Alpine | L | 0.5 | 0.5 | 0.895105 | 232.219168 | 6 |
| Alpine | H | 4 | 0.5 | 0.898295 | 224.3725035 | 5 |
| Alpine | H | 2.5 | 0.5 | 0.906895 | 215.6799301 | 5 |
| Alpine | H | 0.5 | 0.6 | 0.89088 | - | 23 |
| Jackson | H | 3.5 | 0.090909091 | 0.7772409 | 276.9584756 | 6 |
| Jackson | H | 4 | 0.090909091 | 0.7773182 | 296.5895898 | 7 |
| Jackson | H | 3 | 0.090909091 | 0.7773864 | - | 12 |
| Jackson | H | 2 | 0.090909091 | 0.8230182 | - | 12 |
| Jackson | H | 2.5 | 0.090909091 | 0.8234773 | - | 14 |
| Jackson | LQ | 4 | 0.181818182 | 0.7101455 | 251.8684554 | 1 |
| Jackson | L | 3.5 | 0.272727273 | 0.60175 | 258.7359365 | 1 |
| Jackson | L | 4 | 0.272727273 | 0.6024864 | 259.499018 | 1 |
| Jackson | LQ | 3.5 | 0.272727273 | 0.6793591 | 251.7077368 | 2 |
| Jackson | L | 0.5 | 0.272727273 | 0.7016091 | 263.8753876 | 5 |

| | | | | | | |
|---|---|---|---|---|---|---|
| Jackson | L | 1 | 0.272727273 | 0.7086545 | 254.3545529 | 3 |
| Jackson | LQ | 1 | 0.272727273 | 0.8240091 | 251.79814 | 4 |
| Jackson | LQ | 1.5 | 0.272727273 | 0.8318727 | 247.5532945 | 3 |
| Jackson | LQ | 2.5 | 0.272727273 | 0.8378864 | 246.4362601 | 2 |
| Jackson | LQ | 2 | 0.272727273 | 0.8389182 | 244.8357271 | 2 |
| Jackson | L | 2.5 | 0.363636364 | 0.5528318 | 258.4254047 | 2 |
| Jackson | L | 3 | 0.363636364 | 0.5833864 | 258.0480926 | 1 |
| Jackson | L | 2 | 0.363636364 | 0.64355 | 255.1049989 | 2 |
| Jackson | L | 1.5 | 0.363636364 | 0.6852545 | 256.5961569 | 3 |
| Jackson | LQ | 0.5 | 0.363636364 | 0.7837682 | 285.3357205 | 7 |
| Jackson | LQ | 3 | 0.363636364 | 0.7894136 | 248.625469 | 2 |
| Jackson | H | 1.5 | 0.363636364 | 0.8025455 | - | 14 |
| Jackson | H | 1 | 0.454545455 | 0.7811682 | - | 17 |
| Jackson | H | 0.5 | 0.818181818 | 0.6934545 | - | 45 |

**Fig. S1**

a)

**Furthest away locality**



b)

**Expansion source**



**Fig. S2**

**a)**



**b)**



**c)**



**d)**



**Fig. S3**