

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27

DR. QIXIN HE (Orcid ID : 0000-0003-1696-8203)

Article type : Original Article

Title: 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

Qixin He<sup>1,2</sup>, Joyce R. Prado<sup>3</sup>, and L. Lacey Knowles<sup>4</sup>

<sup>2</sup>Department of Ecology and Evolutionary Biology, University of Chicago, Chicago IL, USA 60637

<sup>3</sup>Departamento de Ciências Biológicas, Escola Superior de Agricultura ‘Luiz de Queiroz’, Universidade de São Paulo, Piracicaba, Brazil

<sup>4</sup>Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor MI, USA 41809-1079

<sup>1</sup>Corresponding Author: Qixin He<sup>1</sup>, E-mail: [heqixin@uchicago.edu](mailto:heqixin@uchicago.edu)

Orcid id: 0000-0003-1696-8203

RH: Inferring the origin of range expansion

Contact Information:

Qixin He, [heqixin@uchicago.edu](mailto:heqixin@uchicago.edu)

University of Chicago

208 Erman

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/mec.14380](https://doi.org/10.1111/mec.14380)

This article is protected by copyright. All rights reserved

28 1101 E 57<sup>th</sup> Street  
29 Chicago, IL, 60637

30

31 Word count: Main Text 6613

32 Number of Figures: 5 + 3 SI

33 Number of Tables: 1 + 2 SI

34 Data accessibility: vcf files for Collared pika and SNP used in the analysis were deposited on  
35 Dryad for data archive (DOI: 10.5061/dryad.4s1gg); X-Origin pipeline tutorial, scripts, example  
36 files and input files used in the study are uploaded on GitHub and released under the  
37 DOI: <https://zenodo.org/badge/latestdoi/100994225>.

38

39 ABSTRACT

40

41 Climatic or environmental change is not only driving distributional shifts in species today, but it  
42 has also caused distributions to expand and contract in the past. Inferences about the geographic  
43 locations of past populations, especially regions that served as refugia (i.e., source populations)  
44 and migratory routes are a challenging endeavor. Refugial areas may be evidenced from fossil  
45 records or regions of temporal stability inferred from ecological niche models. Genomic data  
46 offer an alternative and broadly applicable source of information about the locality of refugial  
47 areas, especially relative to fossil data, which are either unavailable or incomplete for most  
48 species. Here we present a pipeline we developed (called X-ORIGIN) for statistically inferring the  
49 geographic origin of range expansion using a spatially explicit coalescent model and an  
50 Approximate Bayesian Computation testing framework. In addition to assessing the probability  
51 of specific latitudinal and longitudinal coordinates of refugial or source populations, such  
52 inferences can also be made accounting for the effects of temporal and spatial environmental  
53 heterogeneity, which may impact migration routes. We demonstrate X-ORIGIN with an analysis of  
54 genomic data collected in the Collared pika that underwent post-glacial expansion across Alaska,  
55 as well as present an assessment of its accuracy under a known model of expansion to validate  
56 the approach.

57

58

60

61 Population expansions leave signatures in the distribution of population genetic variation  
62 across a landscape. This pattern of genetic variation is commonly used for making inferences  
63 about the underlying demographic processes. For example, the decreasing pattern of genetic  
64 diversity along expansion routes has been used to infer the origin of human migrations  
65 (DeGiorgio, Jakobsson, & Rosenberg, 2009; Ramachandran et al., 2005). Similarly, such genetic  
66 signatures have been applied to study post-glacial expansions in other species, as well as their  
67 corresponding geographic refugia during glacial periods of the Pleistocene (reviewed in Hewitt,  
68 2000).

69 However, this approach comes with an inherent issue. Specifically, genetic diversity  
70 patterns (e.g., heterozygosity,  $F_{ST}$ ) can reflect not only signatures from recent distributional  
71 shifts, but also local habitat suitability or long-term geographic isolation (Austerlitz, Jung-  
72 Muller, Godelle, & Gouyon, 1997; Ray, Currat, & Excoffier, 2003). Thus, while the isolation-  
73 by-distance model applies relatively well to species that have a broad habitat, such as human  
74 beings, species with narrower niches tend to track their habitats, displaying a genetic diversity  
75 pattern of isolation-by-barriers or resistance (McRae & Beier, 2007). Therefore, sole reliance on  
76 the gradients of population size/heterozygosity or the principal components without spatial  
77 models is inadequate for making accurate inferences about the ancestral source population or  
78 directions of expansion (François et al., 2010). Due to the rich, yet confounding information  
79 retained in the genetic diversity patterns, most phylogeographic studies infer the location of  
80 hypothesized refugia from the data that are independent of the genomic information (reviewed in  
81 Knowles, 2009). Ecological niche models (ENMs), for instance, could be applied to infer areas  
82 with temporal stability as suitable habitats. In addition, the associated genetic data could then be  
83 used to evaluate the hypothesis that such geographic regions would have served as refugial  
84 source population (e.g., see Carnaval, Hickerson, Haddad, Rodrigues, & Moritz, 2009; Knowles,  
85 Massatti, He, Olson, & Lanier, 2016).

86 Attempts to address the issue of complex historical processes shaping the current genetic  
87 patterns have witnessed the development of spatially-explicit demographic models as well as  
88 spatial genetic indices. Ray, Currat, Berthier, & Excoffier (2005) systematically tested the  
89 likelihood of different geographic locations as human origins by evaluating the goodness-of-fit

90 of  $R_{ST}$  values from different spatial simulations of expansions using the empirical values. Itan,  
91 Powell, Beaumont, Burger, & Thomas (2009) estimated the origin of lactase persistent mutations  
92 in Europe by fitting empirical frequencies of lactase persistent mutations to those from spatial  
93 simulations of the gene expansion along with dairy groups. These pioneer studies demonstrate  
94 the potential of using spatially-explicit models for estimating migration histories. However, these  
95 models do not take temporal changes in habitat suitability into account, which limit their  
96 applicability in flora and fauna that underwent expansions largely driven by climatic oscillations.

97 Spatial genetic indices, on the other hand, are designed to pick up “range expansion”-  
98 specific signatures – that is, the directions of gene flow. By analyzing the allele frequency clines  
99 created by consecutive founder events during the expansion of a population across a landscape,  
100 as captured by a directionality index  $\Psi$ , Peter & Slatkin (2013) demonstrated how information on  
101 the geographic origin and the direction of expansion could be extracted from genomic data  
102 through asymmetrical gene flow. That is, regression between pairwise differences of  $\Psi$  and  
103 geographic distances between populations can be used to *directly* infer the geographic origin of  
104 expansion. However, several aspects of this approach limit its utility in practice. For example,  
105 this method does not account for the heterogeneity in the underlying landscape during the  
106 inference procedure (i.e., assuming a strict isolation-by-distance model).  $\Psi$  may also be biased  
107 toward non-zero values when local population sizes differ substantially (Peter & Slatkin, 2013).  
108 Also, although it is possible to recover a signature of expansion from the magnitude of  $\Psi$ ,  
109 assessing the significance of  $\Psi$ -values, and hence, the confidence of the inferred origin, is not  
110 straightforward.

111 Here, we present a pipeline specifically developed for making statistical inferences about  
112 the geographic origin of range expansion (called X-ORIGIN) that addresses these aforementioned  
113 shortcomings. This pipeline builds upon earlier developments in spatial demographic models  
114 (e.g., Ray et al., 2010) and spatially explicit summary statistics (e.g., Peter & Slatkin, 2013).  
115 Specifically, with the X-ORIGIN we couple the  $\Psi$ -index (Peter & Slatkin, 2013) with a spatially  
116 explicit coalescent model for hypothesis testing in an Approximate Bayesian Computation  
117 (ABC; Beaumont, Zhang, & Balding, 2002) framework. Information based on current and/or  
118 historical habitat suitability can be estimated using ENMs and subsequently incorporated into the  
119 spatially explicit coalescent model (i.e., a modified application of SPLATCHE2; Ray, Currat, Foll,  
120 & Excoffier, 2010). In addition, with the ABC framework, the estimation of the geographic

121 origin of range expansion will not be sensitive to the uncertainties in the underlying demographic  
122 parameters if a wide range of priors of demographic parameters is specified in spatial simulations.  
123 Hereafter, we refer to the geographic origin of range expansion as a parameter,  $\Omega$ . Together, the  
124 significance of expansion and the confidence of a particular geographic location for the ancestral  
125 source population are provided by the X-ORIGIN. As such, the pipeline couples information from  
126 a series of independent analyses (Fig. 1), making X-ORIGIN a useful tool for inferring the  
127 geographic origin of ancestral sources with confidence.

128 It should be noted that there are general procedural parallels with the integrative  
129 distributional, demographic, and coalescent (iDDC) approach for model selection, which also  
130 involves a series of independent analyses (i.e., estimates of habitat suitability, demographic  
131 modeling, and spatially explicit coalescent; He, Edwards, & Knowles, 2013). However, the X-  
132 ORIGIN pipeline differs in that (i) it infers a novel model parameter of interest  $\Omega$  (i.e., the actual  
133 latitudinal and longitudinal coordinates), and (ii) it utilizes information from spatial summary  
134 statistics, specifically, pairwise population measures of  $F_{ST}$  and the directionality index,  $\Psi$  (Peter  
135 & Slatkin, 2013). As such X-ORIGIN is an approach that focuses on the estimation of a specific  
136 parameter of interest –  $\Omega$ , whereas the iDDC is an approach for model selection among a set of  
137 biologically informed demographic hypotheses, the foci of which vary significantly among  
138 studies (e.g., Bemmels, Title, Ortego, & Knowles, 2016; Knowles & Massatti, 2017; Massatti &  
139 Knowles, 2016).

140 Here we describe the approach and test the accuracy of the X-ORIGIN pipeline in inferring  
141  $\Omega$  under a known expansion history (i.e., simulated history; see Fig. 2). Specifically, we model a  
142 history of expansion that involves temporal shifts in the habitat suitability of a landscape (i.e., we  
143 validate the approach by implementing a complex model which cannot be accommodated by any  
144 other currently existing programs). We also demonstrate the utility of the X-ORIGIN with an  
145 analysis of empirical data. Specifically, we analyze the SNP dataset collected in the Collared  
146 pika (*Ochotona collaris*) (i.e., data from Lanier, Massatti, He, Olson, & Knowles, 2015). The  
147 impact of the glaciations is pronounced in small Alaskan mammals (Galbreath, Cook,  
148 Eddingsaas, & DeChaine, 2011; Knowles et al., 2016; Lanier et al., 2015). While previous  
149 analyses in the Collared pikas also suggested that contemporary environmental factors contribute  
150 less to genomic structure than a dynamic history involving the founding of current populations

151 by ancestral source populations (Lanier et al., 2015), the location of putative ancestral source  
152 populations remains unclear.

153

## 154 METHODS

155

### 156 *Statistic Inferences using the X-ORIGIN pipeline*

157

158 The X-ORIGIN pipeline couples information from a series of independent analyses to make  
159 inferences about  $\Omega$ , the geographic location of ancestral source populations, by estimating the  
160 posterior probability of  $\Omega$  under an ABC framework (Fig. 1). Scripts are provided in the X-  
161 ORIGIN pipeline for all the steps involved and a detailed tutorial is provided on GitHub (see  
162 <https://github.com/KnowlesLab/X-ORIGIN>).

163 Briefly, the approach employs a spatially explicit coalescent to generate expected patterns  
164 of genomic variation under a set of priors, including a prior on  $\Omega$  and priors on demographic  
165 parameters of the expansion process (i.e.,  $k$  and  $m$ , the local population sizes and migration rates,  
166 and an ancestral population size,  $N_A$ ). That is, genomic simulations of range expansion are  
167 initiated at different random locations within the geographic range specified by the prior on  $\Omega$   
168 and for different population size and migration rate values. If there is no prior knowledge on  
169 possible geographic origins, all demes on the map used for demographic simulations will be  
170 tested. Otherwise, a prior on  $\Omega$  can be based on the fossil record, or a general candidate region  
171 might be based on the regression between pairwise population differences of  $\Psi$  and geographic  
172 distances (see Peters & Slatkin 2013).

173 To make inferences using X-ORIGIN that consider the effects of spatial and temporal  
174 environmental heterogeneity on the expansion process, X-ORIGIN models the impact of this  
175 environmental heterogeneity on the expansion process. Specifically, heterogeneity in habitat  
176 suitability might be derived from ecological niche models (ENMs) for the present or the past  
177 (Sindato et al., 2016; Waltari et al., 2007), or from information on known barriers (e.g., mountain  
178 ranges, glaciers, and bodies of water; Boehm et al., 2013; Knowles & Massatti, 2017; Waltari &  
179 Hickerson, 2013). These suitability maps are used to inform demographic dynamics associated  
180 with the expansion process by specifying different likely migration events as a function of spatial  
181 and/or temporal environmental heterogeneity. Specifically, the habitat suitability scores for each

182 deme determine local population sizes, thereby influencing the actual number of migrants across  
183 demes per generation. If distributional shifts are induced by climatic changes, then temporal  
184 shifts in habitat suitability can be incorporated into the demographic modeling (i.e., applying  
185 different relative weighting of suitability information from past versus current ENMs to mirror  
186 trends of climatic change; see Brown & Knowles, 2012), given that shifts in connectivity over  
187 time can influence the expansion process, and consequently, the patterns of genetic variation  
188 across the landscape.

189

#### 190 *Programs called up in the X-ORIGIN pipeline*

191

192 In the X-ORIGIN pipeline, demographic and spatially explicit coalescent simulations are  
193 performed in SPLATCHE2 (Ray et al., 2010) in conjunction with a customized script in the X-  
194 ORIGIN pipeline to allow for temporally changing landscapes. Local demographic parameters  
195 (i.e.,  $k$  and  $m$ ) are informed from habitat suitability by scaling these parameters proportionally to  
196 the habitat suitability values of local demes (Fig. 1), which might be temporally dynamic (i.e.,  
197 the habitat suitability for a particular location may change each generation based on shifting  
198 climatic conditions; see Brown & Knowles, 2012). Each generation,  $m$  proportion of the  
199 population migrates out of the local deme; migration occurs to the adjacent four cells (north,  
200 south, west, east). After the exchange of individuals, local demes grow logistically with a rate  $r$ ,  
201 and are regulated by the local carrying capacity (which are also rescaled as a function of the  
202 habitat suitability of a deme);  $r$  can be set to a specific value (e.g., He et al., 2013), and as we do  
203 here ( $r = 1$ ), or it can also be estimated as a parameter. For each time-forward simulation (i.e., a  
204 spatially explicit map of per generation local population sizes and migration events), a series of  
205 corresponding time-backward coalescent genetic simulation are run, with a separate coalescent  
206 simulation generated for each independent locus in the study. The ancestry of an allele will trace  
207 back from the present into ancestral source populations, where the pattern of gene lineage  
208 coalescence across the landscape and the timing of coalescence is defined by the time-forward  
209 local demographic simulations (i.e., the per generation  $k$  and  $m$  parameter values). SNP mutation  
210 models are then used to simulate patterns of genomic variation in SPLATCHE2, where the state of  
211 each SNP is generated across the independent coalescent simulations.

212 To generate patterns of genomic variation to compare with the empirical data, the  
213 simulated datasets are constructed by sampling the same populations (in geographic space), the  
214 same number of individuals, and the same number of SNPs as the empirical scenario. Summary  
215 statistics are calculated for both the empirical and simulated datasets. These include the spatial  
216 summary  $\Psi$  statistics calculated within, between, across all populations, as well as pairwise  
217 population  $F_{ST}$ -values; ARLEQUIN 3.5 (Excoffier & Lischer, 2010) is used to calculate  $F_{ST}$ . Note  
218 that other non-spatial statistics often used in ABC analyses were also considered (e.g.,  $K$ , the  
219 number of haplotypes, and  $H$ , observed heterozygosity). These additional summary statistics are  
220 not used in the analyses presented here because of the lack information they contained under the  
221 expansion scenarios (see Supplementary Fig. S1); however, a user could employ them in X-  
222 ORIGIN if they determine they are relevant to the expansion history under study.

223 The empirical summary statistics are compared to those from the simulated data using  
224 approximate Bayesian computation (ABC), as implemented with ABCESTIMATOR in  
225 ABCTOOLBOX (Wegmann, Leuenberger, Neuenschwander, & Excoffier, 2010). Rather than  
226 conducting ABC analyses directly on the summary statistics, principal components (PCs) are  
227 extracted from all predictor variables to remove the effects of interactions between summary  
228 statistics, as well as to reduce "the curse of dimensionality" (i.e., when too many statistics are  
229 included, the distance between the simulated and empirical values systematically increases,  
230 reducing the accuracy of parameter estimates and making it more difficult to distinguish among  
231 models) (Wegmann & Excoffier, 2010; Wegmann, Leuenberger, & Excoffier, 2009).

232 Five thousand simulations (0.5%) whose transformed summary statistics are closest to  
233 those calculated from the empirical genomic data are retained for estimating the model  
234 parameters (i.e.,  $\Omega$ , the geographic locations of the ancestral source populations, and the  
235 demographic parameters  $k$ ,  $m$ , and  $N_A$ ). In order to jointly estimate the likelihood of a specific  
236 deme as the origin  $\Omega$  (i.e., a specific longitude and latitude), the kernel densities of  $\Omega$  across the  
237 retained simulations was estimated and used as the likelihood. This provides a non-parametric  
238 way of smoothing and estimating the likelihood of the origin based on the limited retained  
239 simulations (i.e., from the 0.5%, or five thousand retained simulations).

240 To check if the inferred model is capable of generating the observed data, the likelihood  
241 of the empirical data given the model is compared with the likelihoods of the retained  
242 simulations. The fraction of simulations that have a smaller likelihood than the empirical data is



243 expressed as a  $P$ -value, with small  $P$ -values indicating that a model is highly unlikely (Wegmann  
244 et al., 2010). Likewise, we conduct standard evaluations of the quality of the inferences from  
245 ABC (e.g., bias in parameter estimates; described below).

#### 246 247 *Performance of the X-ORIGIN pipeline*

248  
249 We tested the pipeline on a simulated scenario (Fig. 2A) to evaluate the performance of the  
250 approach for inferring the geographic location of the source population,  $\Omega$ , under a temporally  
251 changing landscape. Specifically, simulations were conducted on a  $50 \times 50$  deme landscape with  
252 a centrally located geographic barrier that was present in the past but not the present and  
253 expansion proceeding from the upper left deme (Fig. 2A). Simulations were run for 500  
254 generations, in which the barrier persisted for 250 generations. At the end of the simulations, 10  
255 diploid individuals were sampled from 10 demes from across the distributional map. A range of  
256 migration rate, ancestral population size, and carrying capacity values per deme were simulated  
257 to check if the inferred origin is sensitive to particular details of the demographic expansion  
258 process (Table 1).

259 The accuracy of X-ORIGIN was evaluated by measuring the geographic distance between  
260 the actual and inferred geographic location of the source population (i.e., differences in the actual  
261 and inferred latitudinal and longitudinal coordinates). In addition to evaluating the accuracy of  
262 the estimated  $\Omega$  under the model in which expansion proceeded from the upper left deme (Fig.  
263 2), we also tested whether the accuracy of  $\Omega$  varied depending upon the geographic origin of the  
264 expansion. Specifically, we investigated the performance of the model by inspecting the average  
265 error of the inferred  $\Omega$  of 10 pseudo-observed datasets (i.e., PODs from the simulations) in which  
266 the geographic origin of the expansion differed. Specifically,  $\Omega$  was systematically varied so that  
267 each deme across the entire map served as the source of expansion.

268 In addition, the accuracy of X-ORIGIN pipeline is compared with Peter & Slatkin's (2013)  
269 original "time difference of arrival location estimation" (TDOA) approach as well as a modified  
270 TDOA approach, which incorporates spatial heterogeneity in migration patterns (Olave, He, &  
271 Knowles, in prep). Specifically, we calculated the distance between the actual geographic origin  
272 with the one estimated from the TDOA approaches. The TDOA approach identifies the origin of  
273 the expansion by identifying the deme that explains the highest proportion of variation in the

274 correlation of pairwise  $\Psi$  differences and the pairwise differences of geographic distances of the  
275 populations to the potential origin. The modified TDOA approach correlates pairwise  $\Psi$   
276 differences with pairwise resistance differences (McRae & Nürnberger, 2006) in which  
277 heterogeneous landscape is considered (Olave et al., in prep), whereas the original TDOA (Peter  
278 & Slatkin, 2013), assumes migration occurs on a homogeneous landscape (i.e., according to a  
279 random diffusion model). We conducted a cursory examination of the robustness of X-ORIGIN to  
280 model mis-specification as well.

281

### 282 *Demonstration of X-ORIGIN with application to Alaskan Collared pika*

283

284 In addition to details about the ABC analyses (see below), here we briefly describe the empirical  
285 genomic data we analyzed with X-ORIGIN, given that all data used here are from previous  
286 publications and are referenced below. Specifically, we analyze a genomic dataset collected in  
287 the Alaskan Collared pika (for details on library construction and rigorous quality filtering see  
288 Lanier *et al.* 2015). Maps of environmental heterogeneity used in the X-ORIGIN analyses to infer  
289  $\Omega$ , the geographic location of the ancestral source population for the Collared pika, were  
290 generated from ENMs for the present and the last glacial maximum, LGM (see details in  
291 Knowles *et al.* 2016).

292

293 *Genomic dataset.* We analyzed RADseq data for 8 populations; note we excluded the Pika Camp  
294 (Wrangell-St. Elias Mtns; GIS coordinates 61.2170, -138.2670) from our analyses because  
295 previous analyses indicate that it was founded from a separate ancestral refugial source (Lanier *et*  
296 *al.* 2015). Of the 23,493 RADseq loci with at least one biallelic SNP across populations, we  
297 analyzed 6816 loci with one SNP retained per RADseq loci in 50 individuals (i.e., 6-8  
298 individuals per population, with the exception of Jawbone Lake, where  $n = 2$ ); loci in less than  
299 50% of the samples or were not present in more than one individual per population were  
300 excluded. Note that this is an expanded dataset relative to those previously published (i.e., Lanier  
301 *et al.* 2015; Knowles *et al.* 2016) because we recovered more genetic information using ddRAD  
302 aligned to a reference genome for *Ochotona princeps* (American pika; ID: 771).

303 The directionality index  $\Psi$  requires information on the ancestral versus derived states of  
304 SNPs because the statistic is calculated by counting the difference in derived allelic frequencies

305 between pairs of populations (see Eq. 1 in Peter & Slatkin 2013). Ancestral states of independent  
306 biallelic SNPs were determined by aligning the sequences with *Ochotona princeps* (American  
307 pika; ID: 771; <https://www.ncbi.nlm.nih.gov/genome>).

308  
309 *Prior on  $\Omega$ , the geographic locations of the origin of expansion.* The TDOA approach was  
310 conducted to select candidate regions of origin to inform the prior on  $\Omega$  (as opposed to  
311 considering the entire state of Alaska). Specifically, for each potential geographic location as the  
312 site of the ancestral source population (i.e., each deme from the distributional map), linear  
313 regression was performed between pairwise  $\Psi$  differences and the pairwise differences of  
314 geographic distances of the populations to the potential origin. The linear regression was  
315 repeated for each of the different potential geographic origins and the geographic locations with  
316  $R^2$ -values larger than 0.5 were used to specify the prior on the geographic location of the  
317 ancestral source population (regression analyses were conducted using modified scripts from  
318 Peter & Slatkin, 2013, which we provide on KnowlesLab/Github). This generated a target area of  
319 approximately 442,300  $km^2$  (i.e., 1302 demes, with a size of  $18.4 \times 18.4 km^2$  for each deme;  
320 Table 1) to analyze in detail regarding the posterior probability of  $\Omega$ , the geographic location of  
321 the ancestral source population for the set of 8 Collared pika populations collected across its  
322 range (see Lanier et al. 2015 for details).

323  
324 *Estimates of habitat heterogeneity across space and time.* Maps of environmental heterogeneity  
325 for the Collared pika were generated from ENMs (see details in Knowles *et al.* 2016). Briefly,  
326 inferences about differences in habitat suitability across space were made for the present and the  
327 LGM from ENMs based on bioclimatic data for the present and paleoclimatic data from 21 kya.  
328 The models were tested over combinations of regularization parameters from 0.25 to 3 in  
329 intervals of 0.25 and the Linear, Quadratic, Hinge, Product and Threshold features using  
330 SDMTOOLBOX (Brown, 2014). Each model parameter class was replicated 25 times using cross-  
331 validation.

332 In addition, temporal shifts in habitat suitability were represented using differences in the  
333 relative weighting of habitat suitabilities estimated for the present and LGM across time to  
334 reflect climatic trends in the region over the past 21,000 years (Brown & Knowles 2012).  
335 Specifically, the current ENM suitability map was used to represent the present to 5,000 years

336 ago, an intermediate suitability map (i.e., an average suitability between the current and LGM  
337 ENMs) for the time period 5,000 - 11,000 years ago, and the LGM ENM suitability map for  
338 11,000 – 21,000 years ago.

339  
340 *ABC analyses.* Datasets were simulated for 2100 generations (based on a scaling factor of 10 to  
341 reduce the computational requirements; see He et al. 2013) to represent the range expansion from  
342 last glacial maximum. Priors for the local carrying capacity ( $k$ ), ancestral population size ( $N_{ans}$ ),  
343 and migration rates ( $m$ ) are set according to Lanier *et al.* (2015) (see Table 1). Note that a  
344 geographic grid of  $18.4 \times 18.4 \text{ km}^2$  corresponded to a single deme and expansion was modeled  
345 across the Alaskan landscape (i.e., over approximately  $2,197,850 \text{ km}^2$ ).

346 As with tests of the general performance of X-ORIGIN, we compared the estimates of  $\Omega$ ,  
347 the geographic location of the ancestral source of expansion, with results from: 1) the TDOA  
348 method, where heterogeneity in the present landscape is not incorporated (i.e., the geographic  
349 distances separating populations were represented as pairwise Euclidean distances), 2) the  
350 modified TDOA method, where resistance distances based on heterogeneity in the current habitat  
351 suitability is used, and 3) X-ORIGIN, where temporal shifts in the heterogeneity of the landscape  
352 over time are accounted for. To evaluate the accuracy of estimates of  $\Omega$ , five thousand pseudo-  
353 observations were generated from the prior distributions of the parameters. If the estimated  
354 parameters are unbiased, posterior quantiles of the parameters from the pseudo datasets should  
355 be uniformly distributed (Cook, Gelman, & Rubin, 2006; Wegmann et al., 2010). The posterior  
356 quantiles of true parameters for each pseudo run were calculated based on the posterior  
357 distribution of the regression adjusted 5000 simulations closest to the pseudo-observed datasets.

358

## 359 RESULTS

360

### 361 *Performance of the X-ORIGIN pipeline*

362

363 For the example history considered here, which involved a central barrier that was present in the  
364 past, but not the present (i.e., there is both spatial and temporal heterogeneity in habitat  
365 suitabilities) X-ORIGIN gives more accurate inferences of  $\Omega$ , the geographic location of the  
366 source population of the expansion, than the TDOA approach. In fact, the performance of X-

367 ORIGIN was quite good, estimating the most likely origin within 1- 4 demes of the actual origin  
368 (mean  $P$ -value = 0.67) from different starting positions across the map (and hence, differences in  
369 when and where the expansion process interacted with the geographic barrier), except for the  
370 lower left grid of the geographic area (Fig. 3A, C; see Supplementary Fig. S2 for detailed  
371 examples of variation in inferences across PODs for different locations of origins).

372 In contrast, the majority of analyses with the TDOA approach give inferred locations that  
373 differ markedly from the actual area where the expansion originated, irrespective of where on the  
374 map the expansion originates (Fig. 3B). The performance of the TDOA approach was especially  
375 poor (i.e., large discrepancies between the inferred and actual geographic origin of expansion)  
376 when the ancestral source area was near the barrier (Fig. 3D). This variation in accuracy  
377 highlights the importance of explicitly modeling the temporal heterogeneity of landscapes (also  
378 see Wegmann, Currat, & Excoffier, 2006), as it strongly distorts the  $\Psi$  signatures, especially if  
379 the heterogeneity is present in the early stage of the expansion.

380

#### 381 *Inferred geographic origin of expansion in the Alaskan Collared Pika*

382

383 For the set of Collared pika populations studied here, the highest likelihood (marginal  
384 density:  $1.82 \times 10^{-8}$ ;  $P$ -value: 0.996) for the location of the expansion origin,  $\Omega$ , is the Mackenzie  
385 Mountains in Yukon Territory, Canada (Fig. 4). This inference is based on the retained 5000  
386 simulations whose summary statistics were to those of empirical data. The geographic origin of  
387 expansion (i.e., the latitudinal and longitudinal coordinates) was estimated using a two-  
388 dimensional kernel density of the retained simulations, implemented using the kde2d function in  
389 the MASS package of R (Venables & Ripley, 2002).

390 The geographic origin of expansion inferred using X-ORIGIN differed from the TDOA  
391 results (Fig. 4). Moreover, neither the inferred area based on the pairwise  $\Psi$  matrix on a  
392 homogeneous landscape (TDOA-diffusion) nor the one based on a resistance map of the current  
393 landscape suitabilities (TDOA-resistance), are in areas with high likelihoods. That is, simulated  
394 genetic data sets where expansion proceeded from the inferred areas under the TDOA  
395 approaches do not correspond to the observed genetic data (i.e., there is a mismatch between the  
396 empirical summary statistic and those calculated from the simulations).

397 Based on the distances between actual versus inferred origin for each of the different  
398 method, X-ORIGIN outperformed TDOA, although the accuracy of inferred  $\Omega$ -values varied  
399 depending upon the geographic origin of the expansion (Fig. 5). We also note that the accuracy  
400 was generally lower for the heterogeneous landscape inferred for pikas relative to the landscape  
401 used to validate the X-ORIGIN package (Fig. 5 versus Fig. 3). In particular, populations that  
402 originated from the southeast region exhibited the lowest accuracy (i.e., the greatest difference  
403 between the inferred and actual value of  $\Omega$ ). This is most likely due to the lack of samples from  
404 that area, and consequently little information of the direction of asymmetrical gene flow  
405 expected under an expansion model (see Peter & Slatkin, 2013). Nevertheless, comparison of the  
406 accuracy of inferences between X-ORIGIN and TDOA approaches, indicate those from X-ORIGIN  
407 are more accurate for an expansion originating from the Mackenzie Mountain range.  
408 Specifically, analyzing simulated data of expansions from the Mackenzie Mountain range (i.e.,  
409 the PODs from the ABC simulations), the TDOA approaches give estimates that are generally  
410 displaced by 15 to 30 demes from the actual origin of expansion (i.e., a discrepancy of 750 to  
411 1500 km), and curiously these were more inaccurate than inferences with a south-west  
412 geographic origin of expansion (Fig. 5), despite sampling of populations in that region (see  
413 discussion below).

#### 414 415 DISCUSSION

416  
417 Patterns of genetic variation in individuals sampled in the present harbor rich information about  
418 past movements of species. In contrast to those from non-spatial models of population  
419 demography (e.g., changes in population size or admixture proportions; see Hey, 2005; Theis,  
420 Ronco, Indermaur, Salzburger, & Egger, 2014), recent developments have focused on inferences  
421 from spatially explicit approaches. Specifically, departure from equilibrium status of population  
422 movements under a diffuse model, ‘isolation-by-distance’, caused either by range  
423 expansion/contraction history, long distance admixture or habitat heterogeneity is tested through  
424 different approaches. One general approach is to quantify discrepancies between spatial genetic  
425 patterns and the expectations from geographic distances. For example, discrepancies between  
426 population’s positions on a genetic PCA map can be visualized against a map of their  
427 geographical distribution using Procrustes analyses to examine where on a landscape patterns of

428 genetic variation depart from isolation by distance (Knowles et al., 2016; Wang, Zöllner, &  
429 Rosenberg, 2012), or a “geogenetic map” can be used to infer potential long-range admixture  
430 among populations (Bradburd, Ralph, & Coop, 2016). Similarly, disruptions to past movement  
431 might be inferred by relating the effective migration rates to expected genetic dissimilarities for  
432 an interpolated geographical map of barriers or corridors among populations (see Petkova,  
433 Novembre, & Stephens, 2016).

434 Instead of quantifying discrepancies from isolation-by-distance, our approach directly  
435 models expected patterns of genetic variation using spatial genetic indices and makes inferences  
436 about historical movements – specifically, the geographic origin of expansion,  $\Omega$  – under an  
437 ABC framework, while incorporating temporal shifts in habitat suitability over time. This is not  
438 the first approach for directly evaluating genetic variation under models of historical movement.  
439 For example, the spatial genetic indices applied here were developed to directly infer historical  
440 movements based on shifts in the genetic summary statistics across a landscape (Peter & Slatkin,  
441 2013), and spatial-autocorrelation of genetic covariance information has been applied to  
442 distinguish among spatially-explicit demographic scenarios (Alvarado-Serrano & Hickerson,  
443 2016; Bertorelle & Barbujani, 1995; Coop, Witonsky, Rienzo, & Pritchard, 2010). However, our  
444 approach infers and evaluates the parameter  $\Omega$  – the actual latitudinal and longitudinal  
445 coordinates for the origin of an expansion– that is not based on the assumption of a diffusion  
446 model that provides statistical rigorousness and flexible applications for inferences about  
447 historical expansion scenarios. First, we can evaluate the likelihood of different geographic  
448 locations for the origin of a population expansion, accounting for both spatial and temporal  
449 heterogeneity in habitat suitability of the landscape. Second, with the freely available X-ORIGIN  
450 pipeline we developed, users can validate any inference by determining whether the inferred  
451 model is capable of generating data that generally corresponds to the empirical data, which is  
452 equally important as estimating the most likely model for the origin of expansion (i.e., the most  
453 likely location for the origin of expansion may nonetheless be a poor fit to the observed data).  
454 Such attributes are not currently implemented in other methods for inference about expansion  
455 histories (e.g., compare with Ray et al., 2005).

456 Below we discuss how these attributes make X-ORIGIN not only a practical tool, but as our  
457 analyses demonstrate, also one whose performance is better than not accommodating such  
458 dynamic histories. Likewise, we highlight how this pipeline can easily be adapted for a more

459 general inference approach beyond inferring the origin of expansions, especially with the  
460 development of new spatial indices. However, we also note the difference in performance of X-  
461 ORIGIN between a simple demographic history (i.e., the one used to validate the approach) and  
462 the one with more extreme habitat heterogeneity, and caution users of the importance for  
463 validating the accuracy of the inference, which can be implemented in the X-ORIGIN pipeline. We  
464 apply this practice when interpreting the results from the X-ORIGIN analysis of the Collared pikas,  
465 as well as discuss aspects of the data that might contribute to uncertainty in the inferred origin of  
466 expansion, and the importance of corroborative evidence not based on the genetic data itself.

467

#### 468 *Factors impacting the accuracy of inferences about the geographic origin of expansion*

469

470 The  $\Psi$  index directly captures the overall trend of differences in frequencies of derived  
471 polymorphic alleles in populations based on the fact that expanding front of populations are  
472 experiencing serial bottlenecks. Therefore,  $\Psi$  indices are informative as long as current  
473 populations have not yet reached equilibrium. If the majority of the pairwise  $\Psi$  indices are close  
474 to zero in the system (which is not the case for pikas; Supplementary Table 1), the lack of spatial  
475 gradient in the  $\Psi$  indices indicate that either there was not an expansion or a sufficient amount of  
476 time since the expansion has passed such that its genetic signature can no longer be detected by  
477 the  $\Psi$  indices (see also Peter and Slatkin 2013). We tested a scenario in the Pika dataset where  
478 there is no expansion origin to examine the performance of X-ORIGIN. Specifically, we simulated  
479 1000 replicate data sets in which all populations started from their sampling areas to reach  
480 equilibrium states. For these datasets, although  $\Psi$  indices deviate strongly from zero, no origin  
481 can be estimated from TDOA as no positive relationship between pairwise differences of  $\Psi$  and  
482 geographic distances among populations can be established (Supplementary Table 2). Likewise,  
483 with X-ORIGIN, marginal densities of the expansion model are extremely low (on the order of  
484  $10^{-200}$  to  $10^{-12}$  as compared to  $10^{-8}$  for PODs that experienced expansion from a single origin) and  
485  $P$ -values are zero (Supplementary Table 2). Therefore, X-ORIGIN, like TDOA, will not give  
486 misleading results about the potential origin for expansion when no such expansion occurred.

487 Any inference that extracts information on the geographic distribution of genetic variation  
488 requires adequate sampling of populations as well as number of independent SNPs (i.e., at least  
489 more than 1,000 independent SNPs; Peter & Slatkin, 2013; Bradburd et al., 2016). Our results



490 clearly show that inferences become less accurate when sampled populations are located further  
491 from the location where an expansion originated (e.g., see higher error rate at south-east corner  
492 of Fig. 5A). Therefore, researchers should carefully consider the sampling design. In particular,  
493 our results (see also Peter & Slatkin, 2013; Bradburd et al., 2016) suggest that obtaining accurate  
494 inferences that utilize spatial information about the distribution of genetic variation may be  
495 dependent upon which populations are sampled, rather than whether there is sufficient power for  
496 such inferences related to the number of loci analyzed. Although it's beyond the scope of this  
497 study, this general question is something that could be explored using the X-ORIGIN pipeline.

498 Another factor impacting the accuracy of inference relates to model mis-specification.  
499 Specifically, complicated demographic scenarios such as those involving two or more  
500 geographic origins of expansion will give misleading results if not accommodated (see also Peter  
501 & Slatkin 2013). There are a number of ways to accommodate and/or test whether an assumed  
502 expansion from a single source might be violated. For example, clustering algorithms can be run  
503 to delineate populations into different groups with potentially different expansion origins and  
504 validated by a minimum spanning tree built from a matrix of  $\Psi$ -values (Peter & Slatkin, 2013),  
505 followed by separate inferences of  $\Omega$  for each subgroup of populations. Alternatively, competing  
506 explanatory models with multiple origins versus one expansion origin can be analyzed in X-  
507 ORIGIN and compared in a model selection framework. Our results also suggest that any model,  
508 even those that might be more probable than others, should be interpreted with caution if  $\Omega$  is  
509 located in areas with low confidence (based on reference to simulated datasets), or if the most  
510 likely model nevertheless has a low probability of generating data that resembles the empirical  
511 data (i.e., low *P-value*; Wegmann et al. 2010; see He et al. 2013 for details of model validation).

512 Despite positive aspects of X-ORIGIN related to estimating the likelihood of the expansion  
513 origin, and consequently, uncertainty surrounding this inference (e.g., the geographic area  
514 spanned by the 90% highest posterior density of  $\Omega$ ), as well as validating the inference using  
515 PODs (see Fig. 5), one unexplored issue is how errors early in the pipeline might get amplified  
516 and generate misleading results. We did a cursory examination of how such errors might impact  
517 an inferred expansion origin. Specifically, we examined how robust the inferred origin might be  
518 to uncertainties regarding the temporal changes in habitats – in this case, the duration of a  
519 barrier, as in the scenario we used to validate X-ORIGIN (see Fig. 2). When we varied the true  
520 duration of the barrier to simulate data (i.e., simulate data with a barrier that persisted for 200 to

521 300 generations, rather than 250 out of the 500 generations), we observed no difference in the  
522 accuracy of the  $\Omega$  estimation (Supplementary Fig. S3). This shows that the pipeline can be robust  
523 to misspecification of temporal dynamics of a historical scenario (at least for the parameter space  
524 examined here). This clearly should not be interpreted as general evidence of robustness to  
525 model mis-specification. Rather we present it here to show that X-ORIGIN exhibits some  
526 robustness, but also to emphasize that all users can conduct their own investigation to robustness  
527 tailored to the specifics of their application.

528 There are of course other paths for errors that could impact the accuracy of inferences  
529 about  $\Omega$ . For example, we use ENMs to estimate potential suitable areas to inform demographic  
530 models (see Fig. 1). As a consequence, the results from X-ORIGIN could be impacted by poor  
531 ENMs (i.e., validation and best practices of ENMs should be followed). In addition, applying  
532 different transformation of habitat suitabilities into local carrying capacities can affect patterns of  
533 genetic variation (see Brown & Knowles 2012). There are different strategies one might take to  
534 avoid biases that could result from unrealistic assumptions or errors in the upstream steps of the  
535 pipeline (Figure 1). For example, instead of using a fixed suitability score from an ENM model  
536 for each deme, suitability scores between maximum and minimum range inferred for each deme  
537 might be randomly sampled during the simulation process to generate expected patterns of  
538 genetic variation that incorporate some uncertainties in the ENM modeling. This might increase  
539 the number of simulations required for inferring  $\Omega$  to get an unbiased and precise estimate under  
540 an ABC framework, given that accommodating such uncertainties may increase the variance in  
541 observed patterns of genetic variation in simulated datasets. Likewise, different transformations  
542 of habitat suitabilities into local carrying capacities (scaling habitat suitability linearly with local  
543 carrying capacity versus a step function; Brown & Knowles, 2012) could be incorporated as  
544 alternative models to be tested (i.e., treated in a model selection framework, even when the  
545 primary interest is on inferring the origin of expansion,  $\Omega$ ).

546 Although such flexibility in accounting for uncertainty or potential errors in upstream  
547 steps (Fig. 1) is a strength of the X-ORIGIN package we developed, the application of X-ORIGIN  
548 (especially compared with TDOA; Peter & Slatkin, 2013) comes with much more computational  
549 expense. For example, a typical spatially explicit simulation of 2000 generations on a 150 x 150  
550 grid layer and the generation of 1000 SNPs takes more than 7 minutes. Users are advised to  
551 calculate required computational resources before experimenting with the pipeline. This includes

552 reducing the size of the  $\Omega$  prior (e.g., by applying TDOA as a preliminary step for data  
553 inspection, as applied in the Collared pika example).

554

555 *The Mackenzie Mountain region as the most likely origin of expansion in Collared pika*

556

557 As an alpine small mammal, suitable habitats for Collared pika are spatially highly  
558 heterogeneous, but also temporally heterogeneous given that Alaska was directly impacted by  
559 the glacial cycles (Fig. 2B). Previous analyses have suggested a potentially complex  
560 biogeographic history involving expansion from multiple ancestral sources (Knowles et al.,  
561 2016; Lanier et al., 2015; Lanier & Olson, 2009). Limited sampling of populations inhibits  
562 analysis of data subsets to explore such models with X-ORIGIN (i.e., multiple populations are  
563 required to estimate potential sources of expansion), and therefore is beyond the scope of this  
564 manuscript. Nevertheless, it is informative to consider how our inference compares to previous  
565 characterizations for the populations analyzed here.

566 Previous studies that made inferences about the biogeographic and demographic history  
567 of the Collared pika applied analyses that assumed equilibrium status (e.g.,  $F_{ST}$ , STRUCTURE  
568 analyses, estimates of phylogenetic relationships among populations). For example, in an  
569 analysis of the relationship between  $F_{ST}$ -values among populations and the geographic distance  
570 separating them (Lanier et al., 2015), the most north-eastern sampled population Jawbone Lake  
571 (Fig. 4) appeared to be an outlier under the expectation of isolation-by-distance. Based on this  
572 result, and the relative genetic distinctiveness of the Jawbone Lake population and the other two  
573 north-central populations from the Yukon-Tanana Uplands (specifically, the Eagle Summit and  
574 Crescent Creek populations), these populations were analyzed separately and a distinct pattern of  
575 isolation by distance at the regional level was interpreted as possible evidence of different  
576 ancestral source populations (Lanier et al., 2015). However, our analyses here provide a  
577 compelling argument for an alternative explanation. Specifically, the genetic similarities between  
578 Jawbone Lake and the Eagle Summit and Crescent Creek populations (See Fig. 5 in Lanier et al.,  
579 2015) may not reflect a refugial source that was differed from the refugial source of other  
580 sampled populations. Instead, it may reflect their proximity to the geographic origin of expansion  
581 in an ancestral species,  $\Omega$  in the Mackenzie mountains (see Fig. 4), and more specifically, the  
582 similar geographic distance of the populations from the source of expansion. Even though our

583 validation tests indeed show that the degree of reliability about expansion can be considerable  
584 (e.g., differing by as much as 1500 km from the actual expansion origin depending upon where  
585 on the landscape the expansion proceeded from; Fig. 5), the mean error surrounding estimates of  
586  $\Omega$  as a function of the distance from the actual origin is quite low (i.e., less than 5 demes away,  
587 or 250 km) for the geographic region with the highest likelihood of  $\Omega$  (Fig. 4). Interestingly,  
588 Procrustes analyses in the Collared pikas, as well as other co-distributed alpine mammals,  
589 suggest a stronger deviation along the longitudinal axis between genetic variation and geography  
590 (i.e., genetic similarities more centrally located than the geographic space occupied by the  
591 populations; Knowles et al., 2016). Our analysis supported this deviation as a result of an  
592 expansion history along this axis, offering an alternative interpretation to the hypothesis of a  
593 centrally located refugium.

594 Lastly, ENMs for the LGM are not inconsistent with our estimate (Fig. 2B). However, if  
595 we consider information from the ENMs by themselves, the region of high habitat suitability  
596 encompasses a broad area that does not offer much detail about the potential location of ancestral  
597 populations. This even includes a potential north-western source population (Fig. 2B), even  
598 though former genetic (Knowles et al., 2016; Lanier et al., 2015) and fossil studies (Gunderson,  
599 Jacobsen, & Olson, 2009; Lanier & Olson, 2013) suggest the lack of support for such a putative  
600 ancestral source (e.g., in the Brooks Range). Both X-ORIGIN and TDOA analyses reinforce that  
601 despite projections from the ENM for the LGM, this region does not appear to be a likely  
602 candidate as an ancestral source of expansion.

603

604

## 605 CONCLUSIONS

606 Our results show that failing to consider the impact of spatial and temporal heterogeneity  
607 on the expansion process can lead to much less accurate inferences (Fig. 3A compared with 3B,  
608 and Fig. 5A compared with 5B). Furthermore, there are also ways to minimize potential errors  
609 when inferring the origin of expansion. For example, in our simulations, we place a broad prior  
610 on parameters that are not targets of interest, but may influence estimates of  $\Omega$  (e.g., ancestral  
611 population, carrying capacity; see Table 1), thereby accounting for uncertainty about the  
612 demography of the expansion process. Moreover, the summary statistics used in the inference  
613 procedure (i.e.,  $\Psi$  and  $F_{ST}$ -values) are not sensitive to the absolute effective population sizes, but

614 rather the ratio of size differences between population pairs. Lastly, despite the lower accuracy of  
615 inferences for complicated scenarios, as with the analysis of the Collared pika, relative to simple  
616 expansion scenarios (Fig. 3, 5), accounting for the effects of spatial and temporal heterogeneity  
617 is generally more accurate than applying an oversimplified model if the goal is to infer the  
618 geographic location of an expansion's origin (Fig. 3). Therefore, we argue that the caveats and  
619 concerns associated with inferring the origin of expansion do not nullify the utility of spatially  
620 and temporally explicit models, such as those applied here in the new X-ORIGIN pipeline. In  
621 particular, we show that it is incorrect to assume that environmental heterogeneity (whether  
622 temporal or spatial) will not impact inferred origins of expansion, and that despite the caveats we  
623 highlight with X-ORIGIN, they are less problematic than many implicit assumptions made in  
624 approaches that ignore geographic and temporal constraints on population movements or  
625 population size fluctuations (see Knowles & Alvarado-Serrano 2010). Moreover, the reliability  
626 of any inference about the origin of expansion under the more complex models implemented in  
627 the X-ORIGIN pipeline can be (and should be) rigorously explored using validation procedures.

628

629

630

#### 631 ACKNOWLEDGEMENTS

632

633 Specimens and genomic data analyzed here were collected as part of past collaborative work  
634 with Hayley Lanier and Link Olson, for which we are grateful for their helpful discussions on  
635 small Alaskan mammals, and Collared pikas in particular. We appreciate the comments from two  
636 anonymous reviewers and Benjamin Peter that helped improve earlier versions of the  
637 manuscript. The authors declare no conflict of interest.

638

#### 639 REFERENCES

640 Alvarado-Serrano, D. F., & Hickerson, M. J. (2016). Spatially explicit summary statistics for  
641 historical population genetic inference. *Methods in Ecology and Evolution*, 7(4), 418–  
642 427. doi:10.1111/2041-210X.12489

- 643 Austerlitz, F., Jung-Muller, B., Godelle, B., & Gouyon, P.-H. (1997). Evolution of Coalescence  
644 Times, Genetic Diversity and Structure during Colonization. *Theoretical Population*  
645 *Biology*, 51(2), 148–164. doi:10.1006/tpbi.1997.1302
- 646 Beaumont, M. A., Zhang, W., & Balding, D. J. (2002). Approximate Bayesian Computation in  
647 Population Genetics. *Genetics*, 162(4), 2025–2035.
- 648 Bemmels, J. B., Title, P. O., Ortego, J., & Knowles, L. L. (2016). Tests of species-specific  
649 models reveal the importance of drought in postglacial range shifts of a Mediterranean-  
650 climate tree: insights from integrative distributional, demographic and coalescent  
651 modelling and ABC model selection. *Molecular Ecology*, 25(19), 4889–4906.  
652 doi:10.1111/mec.13804
- 653 Bertorelle, G., & Barbujani, G. (1995). Analysis of DNA diversity by spatial autocorrelation.  
654 *Genetics*, 140(2), 811–819.
- 655 Boehm, J. T., Woodall, L., Teske, P. R., Lourie, S. A., Baldwin, C., Waldman, J., & Hickerson,  
656 M. (2013). Marine dispersal and barriers drive Atlantic seahorse diversification. *Journal*  
657 *of Biogeography*, 40(10), 1839–1849. doi:10.1111/jbi.12127
- 658 Bradburd, G. S., Ralph, P. L., & Coop, G. M. (2016). A Spatial Framework for Understanding  
659 Population Structure and Admixture. *PLOS Genetics*, 12(1), e1005703.  
660 doi:10.1371/journal.pgen.1005703
- 661 Brown, J. L. (2014). SDMtoolbox: a python-based GIS toolkit for landscape genetic,  
662 biogeographic and species distribution model analyses. *Methods in Ecology and*  
663 *Evolution*, 5(7), 694–700. doi:10.1111/2041-210X.12200
- 664 Brown, J. L., & Knowles, L. L. (2012). Spatially explicit models of dynamic histories:  
665 examination of the genetic consequences of Pleistocene glaciation and recent climate  
666 change on the American Pika. *Molecular Ecology*, 21(15), 3757–3775.  
667 doi:10.1111/j.1365-294X.2012.05640.x
- 668 Carnaval, A. C., Hickerson, M. J., Haddad, C. F. B., Rodrigues, M. T., & Moritz, C. (2009).  
669 Stability Predicts Genetic Diversity in the Brazilian Atlantic Forest Hotspot. *Science*,  
670 323(5915), 785–789. doi:10.1126/science.1166955
- 671 Cook, S. R., Gelman, A., & Rubin, D. B. (2006). Validation of Software for Bayesian Models  
672 Using Posterior Quantiles. *Journal of Computational and Graphical Statistics*, 15(3),  
673 675–692. doi:10.1198/106186006X136976

674 Coop, G., Witonsky, D., Rienzo, A. D., & Pritchard, J. K. (2010). Using Environmental  
675 Correlations to Identify Loci Underlying Local Adaptation. *Genetics*, *185*(4), 1411–1423.  
676 doi:10.1534/genetics.110.114819

677 DeGiorgio, M., Jakobsson, M., & Rosenberg, N. A. (2009). Explaining worldwide patterns of  
678 human genetic variation using a coalescent-based serial founder model of migration  
679 outward from Africa. *Proceedings of the National Academy of Sciences*, *106*(38), 16057–  
680 16062. doi:10.1073/pnas.0903341106

681 Excoffier, L., & Lischer, H. E. L. (2010). Arlequin suite ver 3.5: a new series of programs to  
682 perform population genetics analyses under Linux and Windows. *Molecular Ecology*  
683 *Resources*, *10*(3), 564–567. doi:10.1111/j.1755-0998.2010.02847.x

684 François, O., Currat, M., Ray, N., Han, E., Excoffier, L., & Novembre, J. (2010). Principal  
685 Component Analysis under Population Genetic Models of Range Expansion and  
686 Admixture. *Molecular Biology and Evolution*, *27*(6), 1257–1268.  
687 doi:10.1093/molbev/msq010

688 Galbreath, K. E., Cook, J. A., Eddingsaas, A. A., & DeChaine, E. G. (2011). Diversity and  
689 Demography in Beringia: Multilocus Tests of Paleodistribution Models Reveal the  
690 Complex History of Arctic Ground Squirrels. *Evolution*, *65*(7), 1879–1896.  
691 doi:10.1111/j.1558-5646.2011.01287.x

692 Gunderson, A. M., Jacobsen, B. K., & Olson, L. E. (2009). Revised Distribution of the Alaska  
693 Marmot, *Marmota broweri*, and Confirmation of Parapatry with Hoary Marmots. *Journal*  
694 *of Mammalogy*, *90*(4), 859–869. doi:10.1644/08-MAMM-A-253.1

695 He, Q., Edwards, D. L., & Knowles, L. L. (2013). Integrative Testing of How Environments  
696 from the Past to the Present Shape Genetic Structure Across Landscapes. *Evolution*,  
697 *67*(12), 3386–3402. doi:10.1111/evo.12159

698 Hewitt, G. (2000). The genetic legacy of the Quaternary ice ages. *Nature*, *405*(6789), 907.

699 Hey, J. (2005). On the Number of New World Founders: A Population Genetic Portrait of the  
700 Peopling of the Americas. *PLOS Biology*, *3*(6), e193. doi:10.1371/journal.pbio.0030193

701 Itan, Y., Powell, A., Beaumont, M. A., Burger, J., & Thomas, M. G. (2009). The Origins of  
702 Lactase Persistence in Europe. *PLOS Computational Biology*, *5*(8), e1000491.  
703 doi:10.1371/journal.pcbi.1000491

704 Knowles, L. L., & Alvarado-Serrano, D. F. (2010). Exploring the population genetic  
705 consequences of the colonization process with spatio-temporally explicit models: insights  
706 from coupled ecological, demographic and genetic models in montane grasshoppers.  
707 *Molecular Ecology*, 19(17), 3727–3745. doi:10.1111/j.1365-294X.2010.04702.x

708 Knowles, L. L., & Massatti, R. (2017). Distributional shifts – not geographic isolation – as a  
709 probable driver of montane species divergence. *Ecography*, n/a-n/a.  
710 doi:10.1111/ecog.02893

711 Knowles, L. L., Massatti, R., He, Q., Olson, L. E., & Lanier, H. C. (2016). Quantifying the  
712 similarity between genes and geography across Alaska’s alpine small mammals. *Journal*  
713 *of Biogeography*, 43(7), 1464–1476. doi:10.1111/jbi.12728

714 Lanier, H. C., Massatti, R., He, Q., Olson, L. E., & Knowles, L. L. (2015). Colonization from  
715 divergent ancestors: glaciation signatures on contemporary patterns of genomic variation  
716 in Collared Pikas (*Ochotona collaris*). *Molecular Ecology*, 24(14), 3688–3705.  
717 doi:10.1111/mec.13270

718 Lanier, H. C., & Olson, L. E. (2009). Inferring divergence times within pikas (*Ochotona* spp.)  
719 using mtDNA and relaxed molecular dating techniques. *Molecular Phylogenetics and*  
720 *Evolution*, 53(1), 1–12. doi:10.1016/j.ympev.2009.05.035

721 Lanier, H. C., & Olson, L. E. (2013). Deep barriers, shallow divergences: reduced  
722 phylogeographical structure in the collared pika (Mammalia: Lagomorpha: *Ochotona*  
723 *collaris*). *Journal of Biogeography*, 40(3), 466–478. doi:10.1111/jbi.12035

724 Massatti, R., & Knowles, L. L. (2016). Contrasting support for alternative models of genomic  
725 variation based on microhabitat preference: species-specific effects of climate change in  
726 alpine sedges. *Molecular Ecology*, 25(16), 3974–3986. doi:10.1111/mec.13735

727 McRae, B. H., & Beier, P. (2007). Circuit theory predicts gene flow in plant and animal  
728 populations. *Proceedings of the National Academy of Sciences*, 104(50), 19885–19890.  
729 doi:10.1073/pnas.0706568104

730 McRae, B. H., & Nürnberger, B. (2006). Isolation by resistance. *Evolution*, 60(8), 1551–1561.  
731 doi:10.1554/05-321.1

732 Olave, M., He, Q., & Knowles, L. L. (in prep). Evidence for shared refugia based on allele  
733 frequency asymmetries of genomic data among five alpine Alaskan small mammal  
734 species.



- 735 Peter, B. M., & Slatkin, M. (2013). Detecting Range Expansions from Genetic Data. *Evolution*,  
736 67(11), 3274–3289. doi:10.1111/evo.12202
- 737 Petkova, D., Novembre, J., & Stephens, M. (2016). Visualizing spatial population structure with  
738 estimated effective migration surfaces. *Nature Genetics*, 48(1), 94–100.  
739 doi:10.1038/ng.3464
- 740 Ramachandran, S., Deshpande, O., Roseman, C. C., Rosenberg, N. A., Feldman, M. W., &  
741 Cavalli-Sforza, L. L. (2005). Support from the relationship of genetic and geographic  
742 distance in human populations for a serial founder effect originating in Africa.  
743 *Proceedings of the National Academy of Sciences of the United States of America*,  
744 102(44), 15942–15947. doi:10.1073/pnas.0507611102
- 745 Ray, N., Currat, M., Berthier, P., & Excoffier, L. (2005). Recovering the geographic origin of  
746 early modern humans by realistic and spatially explicit simulations. *Genome Research*,  
747 15(8), 1161–1167. doi:10.1101/gr.3708505
- 748 Ray, N., Currat, M., & Excoffier, L. (2003). Intra-Deme Molecular Diversity in Spatially  
749 Expanding Populations. *Molecular Biology and Evolution*, 20(1), 76–86.  
750 doi:10.1093/molbev/msg009
- 751 Ray, N., Currat, M., Foll, M., & Excoffier, L. (2010). SPLATCHE2: a spatially explicit  
752 simulation framework for complex demography, genetic admixture and recombination.  
753 *Bioinformatics*, 26(23), 2993–2994. doi:10.1093/bioinformatics/btq579
- 754 Sindato, C., Stevens, K. B., Karimuribo, E. D., Mboera, L. E. G., Paweska, J. T., & Pfeiffer, D.  
755 U. (2016). Spatial Heterogeneity of Habitat Suitability for Rift Valley Fever Occurrence  
756 in Tanzania: An Ecological Niche Modelling Approach. *PLOS Neglected Tropical*  
757 *Diseases*, 10(9), e0005002. doi:10.1371/journal.pntd.0005002
- 758 Theis, A., Ronco, F., Indermaur, A., Salzburger, W., & Egger, B. (2014). Adaptive divergence  
759 between lake and stream populations of an East African cichlid fish. *Molecular Ecology*,  
760 23(21), 5304–5322. doi:10.1111/mec.12939
- 761 Venables, W. N., & Ripley, B. D. (2002). *Modern Applied Statistics with S* (Fourth). New York,  
762 NY: Springer.
- 763 Waltari, E., & Hickerson, M. J. (2013). Late Pleistocene species distribution modelling of North  
764 Atlantic intertidal invertebrates. *Journal of Biogeography*, 40(2), 249–260.  
765 doi:10.1111/j.1365-2699.2012.02782.x

766 Waltari, E., Hijmans, R. J., Peterson, A. T., Nyári, Á. S., Perkins, S. L., & Guralnick, R. P.  
767 (2007). Locating Pleistocene Refugia: Comparing Phylogeographic and Ecological Niche  
768 Model Predictions. *PLOS ONE*, 2(7), e563. doi:10.1371/journal.pone.0000563

769 Wang, C., Zöllner, S., & Rosenberg, N. A. (2012). A Quantitative Comparison of the Similarity  
770 between Genes and Geography in Worldwide Human Populations. *PLOS Genetics*, 8(8),  
771 e1002886. doi:10.1371/journal.pgen.1002886

772 Wegmann, D., Currat, M., & Excoffier, L. (2006). Molecular Diversity After a Range Expansion  
773 in Heterogeneous Environments. *Genetics*, 174(4), 2009–2020.  
774 doi:10.1534/genetics.106.062851

775 Wegmann, D., & Excoffier, L. (2010). Bayesian Inference of the Demographic History of  
776 Chimpanzees. *Molecular Biology and Evolution*, 27(6), 1425–1435.  
777 doi:10.1093/molbev/msq028

778 Wegmann, D., Leuenberger, C., & Excoffier, L. (2009). Efficient Approximate Bayesian  
779 Computation Coupled With Markov Chain Monte Carlo Without Likelihood. *Genetics*,  
780 182(4), 1207–1218. doi:10.1534/genetics.109.102509

781 Wegmann, D., Leuenberger, C., Neuenschwander, S., & Excoffier, L. (2010). ABCtoolbox: a  
782 versatile toolkit for approximate Bayesian computations. *BMC Bioinformatics*, 11, 116.  
783 doi:10.1186/1471-2105-11-116

784  
785 TABLE

786 Table 1. Prior ranges for demographic and genetic parameters used in the demographic  
787 simulations.

Parameters	Description	Prior Ranges	Distribution
<i>m</i>	migration rate between demes	( $10^{-3.6}$ , $10^{-2}$ )	log-uniform
<i>N<sub>ans</sub></i>	ancestral population size before expansion	(36,880, 508,318)	uniform
<i>K</i>	carry capacity per deme	( $10^{3.3}$ , $10^{4.6}$ )	log-uniform
<b>Lat</b>	latitude range of origin	(1,073,893, 1,850,478)	uniform
<b>Long</b>	longitude range of origin	(616,487, 899,496)	uniform

788

789

790 FIGURE LEGENDS

791

792 Figure 1. The required data inputs (shown in boxes) and workflow of the X-ORIGIN pipeline are  
793 highlighted in the schematic. Specifically, to infer the geographic location from which an  
794 expansion originates,  $\Omega$  (i.e., the actual latitudinal and longitudinal coordinates of the ancestral  
795 source population), a habitat suitability map, candidate regions of  $\Omega$ , and priors for demographic  
796 parameters are required. To consider how habitat heterogeneity might impact the range  
797 expansion process, the habitat suitability map can be informed by spatial (as well as temporal)  
798 variation in suitability (e.g., from ENMs based on contemporary bioclimatic variables, or  
799 paleoclimatic variables; see He et al. 2013). Otherwise, the expansion process can be modeled as  
800 a diffusion process (i.e., equal habitat suitability across space and time). Likewise, users have the  
801 option of either entering candidate regions of  $\Omega$  (e.g., a region identified by the regression  
802 approach of Peter and Slatkin 2013; as discussed in the text), or the entire map area can be  
803 evaluated during the inference procedure. The pipeline calls up different software packages for  
804 downstream generation of simulations and estimation of the expansion origin, candidate regions  
805 of  $\Omega$ . Specifically, spatially explicit coalescent simulations are used to generate expected patterns  
806 of genetic variation under a demographic model the expansion process (either informed or not by  
807 spatial and temporal heterogeneity of the landscape) using a modified version of the program  
808 SPLATCHE2 (Ray et al. 2010). Summary statistics are calculated from each simulated data set  
809 using R script that are incorporated in the pipeline, which are compared with those calculated for  
810 empirical data to inform the posterior distribution of  $\Omega$  using ABC. Note that all steps can be  
811 performed seamlessly in X-ORIGIN, which has a wrapper for connecting all the steps in R or  
812 python scripts. Scripts for the pipeline are shown in grey shaded boxes, while external programs  
813 called in the pipeline are shown without boxes.

814

815 Figure 2. Simulated scenario used to evaluate the performance of the X-ORIGIN pipeline for  
816 inferring the geographic origin of a range expansion. In the simulated scenario, A) expansion  
817 proceeded from the lower-left corner of the map (shown as the red dotted area) across a  
818 homogeneous landscape with a centrally located geographic barrier during the first 250  
819 generations, but not the last 250 generations (i.e., there is spatial and temporal habitat  
820 heterogeneity, where the area of the barrier has zero suitability). Due to the symmetry of the  
821 landscape, we varied the origin of expansion in the simulations within the red dotted area instead

822 of the whole map. Circles mark populations that are sampled and for which summary statistics  
823 are calculated from multiple individuals. B) An empirical application of X-ORIGIN in the  
824 Collared pika in which habitat suitability varied spatially and temporally across the Alaskan  
825 landscape. Ecological niche models were used to estimate habitat suitabilities for the present and  
826 past (i.e., the LGM) using climatic data (see Lanier et al. 2015 for details about ENMs).  
827 Specifically, the demographic expansion process proceeded across a temporally and spatially  
828 heterogeneous landscape, in which the habitat suitabilities from an ENM estimated for the LGM  
829 was used to inform the first 5,000 years of the simulated demographic expansion, followed by  
830 6,000 simulated years of expansion across an intermediate surface (i.e., a map with average  
831 habitat suitability scores between those from the ENM for the present and LGM), and then  
832 10,000 years of expansion with the habitat suitabilities from an ENM based on current climatic  
833 conditions.

834  
835 Figure 3. Distribution of the mean errors in the estimated  $\Omega$  across the map (i.e., for different  
836 geographic locations for the origin of expansion) under the simulated scenario (see Fig. 2A, the  
837 red dotted area) using (A) X-ORIGIN versus (B) the TDOA approach. Color of each deme shows  
838 the accuracy of origin estimation if the expansion starts from that particular deme, which is  
839 measured by the distance between its inferred origin,  $\Omega$ , and the actual origin, averaged across 10  
840 simulations. Also shown are the histograms of accuracy across all 5000 instances (C) from X-  
841 ORIGIN versus the TDOA approach. Distances are in the units of the number of demes from the  
842 actual origin.

843  
844 Figure 4. Estimates of the origin of expansion,  $\Omega$ , inferred in the Collared pika using X-ORIGIN  
845 compared with the TDOA approach. The deme with the highest likelihood inferred from X-  
846 ORIGIN is marked with a black “X”, whereas the location of origins estimated using TDOA  
847 methods are marked with crosses. The heat map shows differences in the probability density  
848 estimates of different demes across the map being the origin of expansion, as estimated in X-  
849 ORIGIN, with the greener shades representing higher probabilities; the shaded square area  
850 represents the prior area for  $\Omega$ , a region identified by the regression approach of Peter and  
851 Slatkin (2013). Each deme in the map has equal relative size (i.e., the map is projected using the  
852 North American Datum – NAD83 – readjustment of the global positioning system that accounts

853 for the earth's curvature) and population localities of sequenced individuals are marked by grey  
854 circles.

855

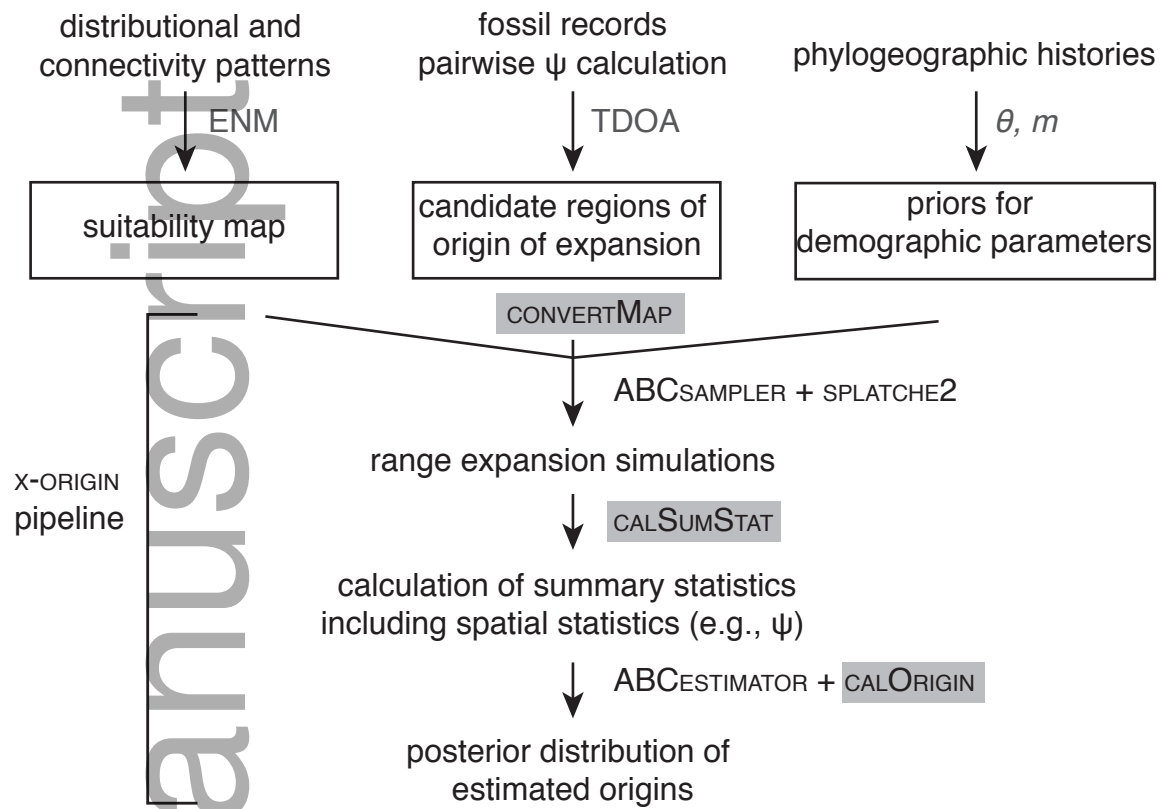
856 Figure 5. Distribution of mean errors in the estimated  $\Omega$  across the map (i.e., for different  
857 geographic locations for the origin of expansion) for pseudo-observations in the Pika simulations  
858 (see Fig. 2B) using (A) X-ORIGIN versus (B) the TDOA approach. 5000 pseudo-observations are  
859 generated and color of each deme shows the accuracy of origin estimation if the expansion starts  
860 from the particular deme, which is measured by the average distance between its inferred origins  
861 and the actual origin. White area on the map contains demes where not all populations can be  
862 colonized if the expansion starts from there. Also shown are the histograms of accuracy across  
863 all 5000 instances (C) from X-ORIGIN versus the TDOA approach. Distances are in the units of  
864 the number of demes from the actual origin and each deme is 18.4km in length.

Author Manuscript

TABLE

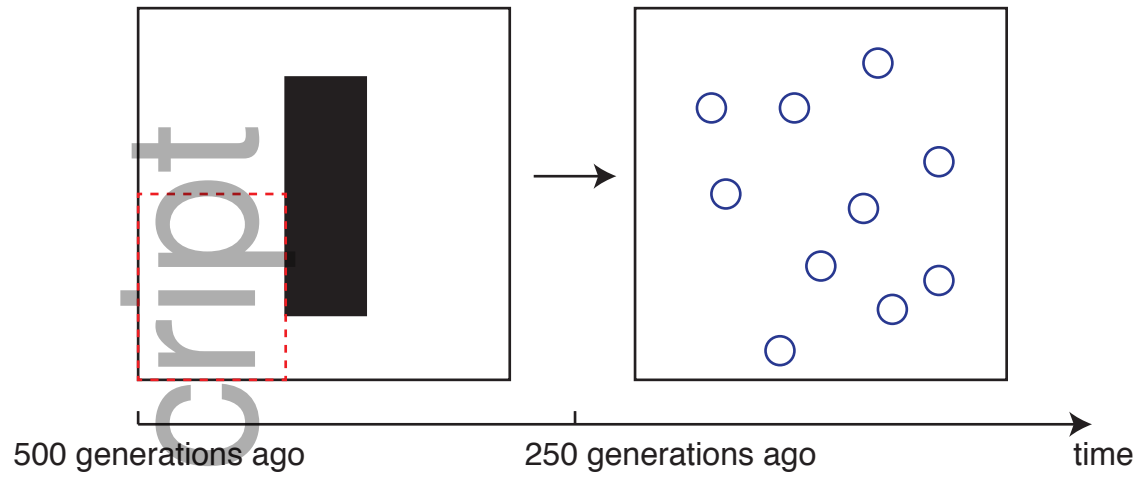
Table 1. Prior ranges for demographic and genetic parameters used in the demographic simulations.

<b>Parameters</b>	<b>Description</b>	<b>Prior Ranges</b>	<b>Distribution</b>
<b>m</b>	migration rate between demes	$(10^{-3.6}, 10^{-2})$	log-uniform
<b>N<sub>ans</sub></b>	ancestral population size before expansion	(36,880, 508,318)	uniform
<b>K</b>	carry capacity per deme	$(10^{3.3}, 10^{4.6})$	log-uniform
<b>Lat</b>	latitude range of origin	(1,073,893, 1,850,478)	uniform
<b>Long</b>	longitude range of origin	(616,487, 899,496)	uniform



A

50 x 50 cells

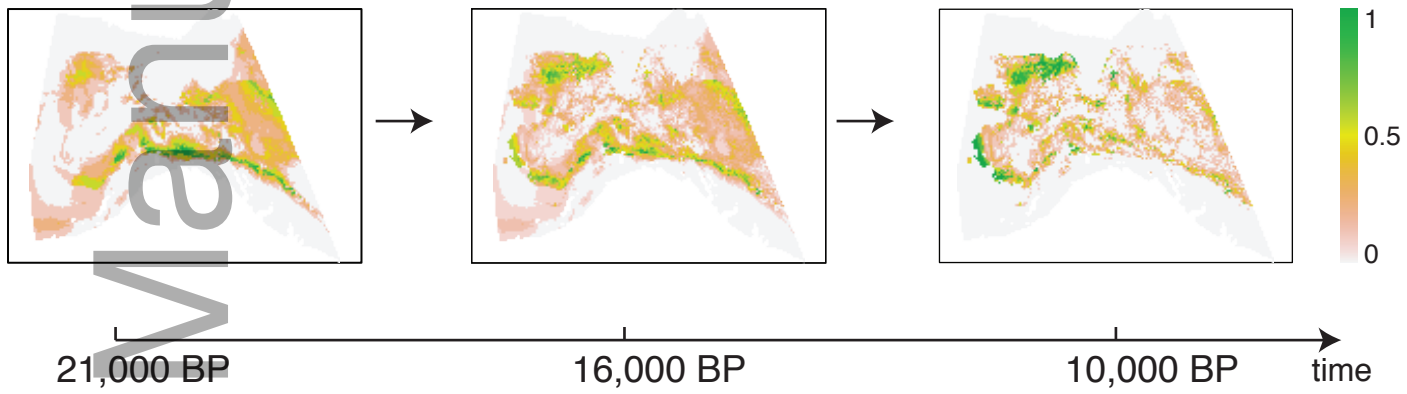


B

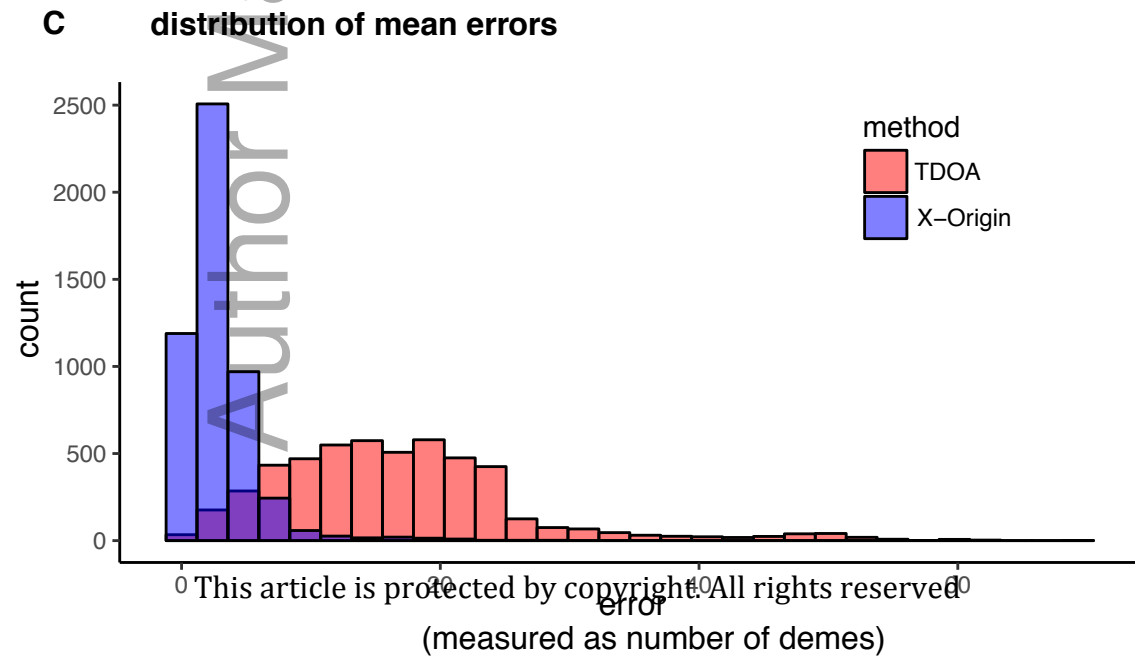
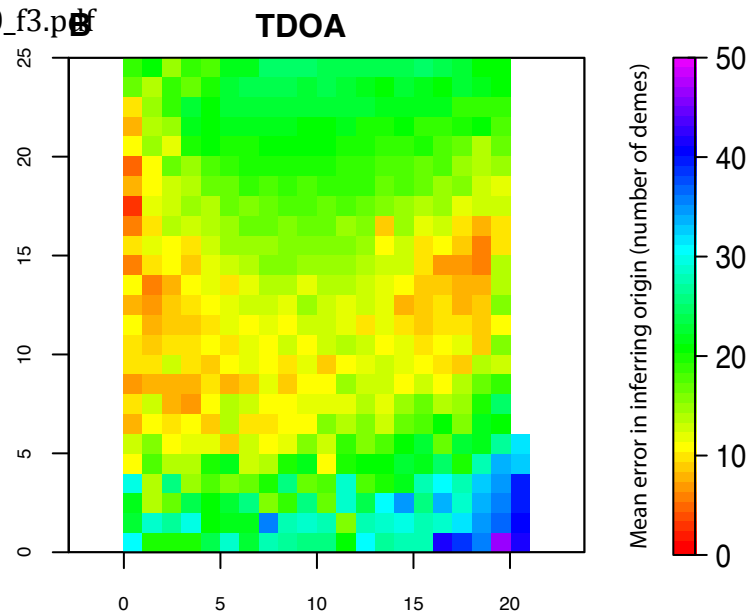
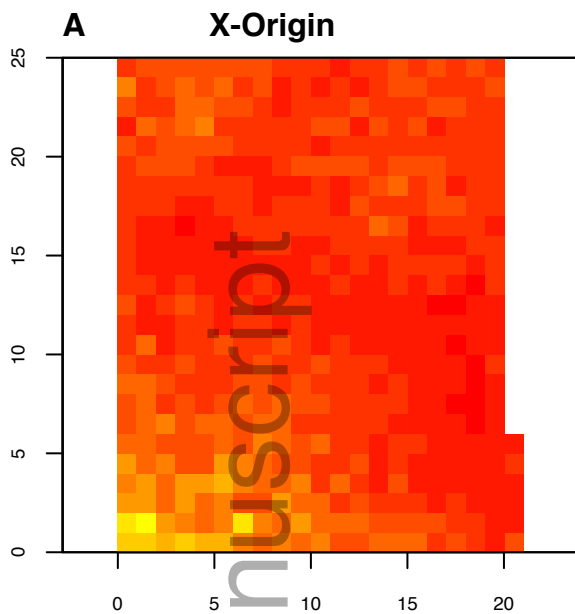
ENM of pika in LGM

Intermediate

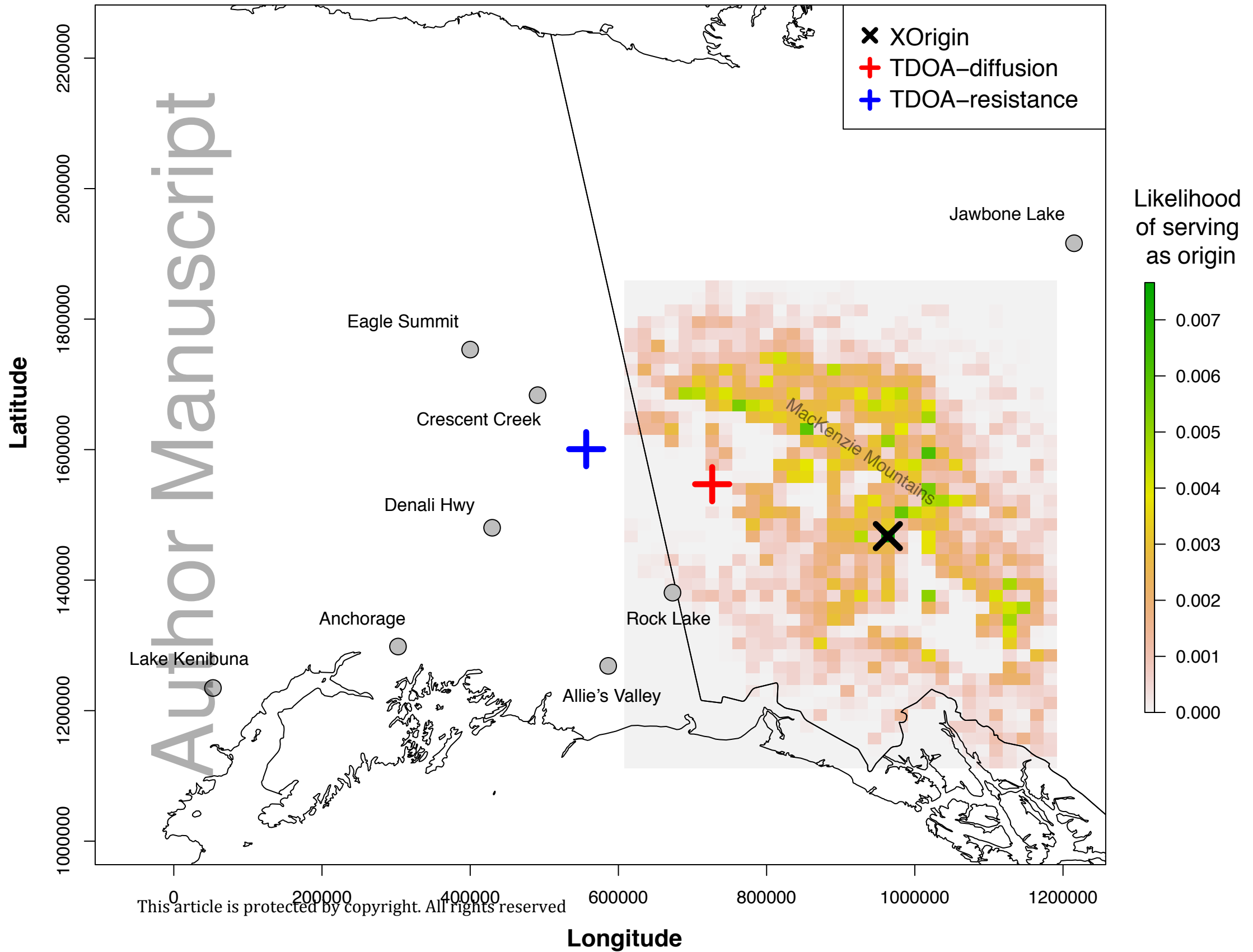
Current

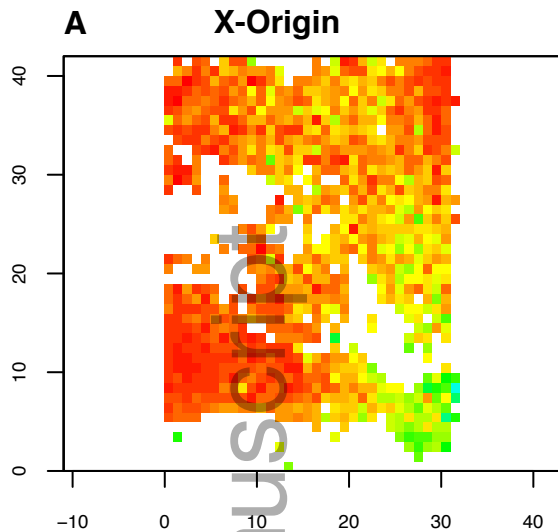




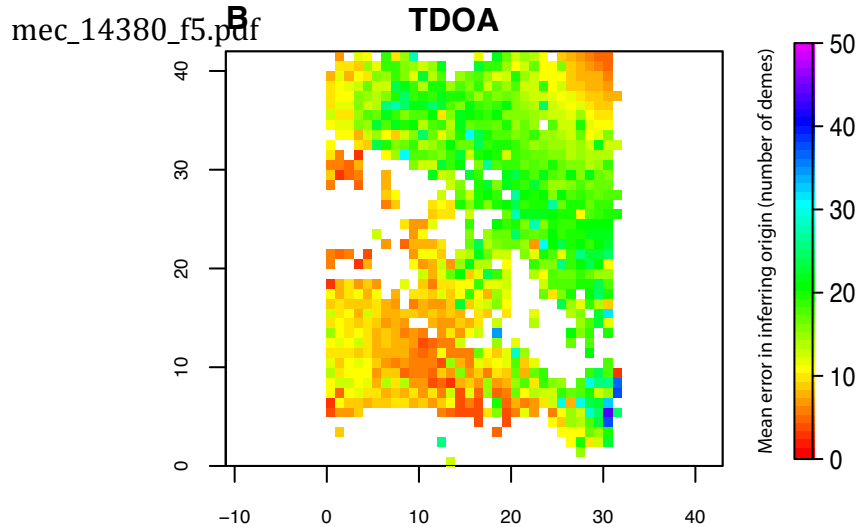


## Comparisons of the origin inferences





**C** distribution of mean errors



**D** distribution of mean errors

