

1 *slendr*: a framework for spatio-temporal 2 population genomic simulations on 3 geographic landscapes

4 Martin Petr¹, Benjamin C. Haller², Peter L. Ralph³, Fernando Racimo¹

5 1. Lundbeck Foundation GeoGenetics Centre, Globe Institute, University of Copenhagen, Denmark

6 2. Department of Computational Biology, Cornell University, Ithaca, NY, USA

7 3. Institute of Ecology and Evolution, University of Oregon, Eugene, OR, USA

8 Abstract

9 One of the goals of population genetics is to understand how evolutionary forces shape patterns
10 of genetic variation over time. However, because populations evolve across both time and space,
11 most evolutionary processes also have an important spatial component, acting through
12 phenomena such as isolation by distance, local mate choice, or uneven distribution of resources.
13 This spatial dimension is often neglected, partly due to the lack of tools specifically designed for
14 building and evaluating complex spatio-temporal population genetic models. To address this
15 methodological gap, we present a new framework for simulating spatially-explicit genomic data,
16 implemented in a new R package called *slendr* (www.slendr.net), which leverages a SLiM
17 simulation back-end script bundled with the package. With this framework, the users can
18 programmatically and visually encode spatial population ranges and their temporal dynamics (i.e.,
19 population displacements, expansions, and contractions) either on real Earth landscapes or on
20 abstract custom maps, and schedule splits and gene-flow events between populations using a
21 straightforward declarative language. Additionally, *slendr* can simulate data from traditional, non-
22 spatial models, either with SLiM or using an alternative built-in coalescent *msprime* back end.
23 Together with its R-idiomatic interface to the *tskit* library for tree-sequence processing and
24 analysis, *slendr* opens up the possibility of performing efficient, reproducible simulations of spatio-
25 temporal genomic data entirely within the R environment, leveraging its wealth of libraries for
26 geospatial data analysis, statistics, and visualization. Here, we present the design of the *slendr* R
27 package and demonstrate its features on several practical example workflows.

28 Introduction

29 Most evolutionary processes in nature have a spatial dimension. Indeed, since its beginnings, the
30 field of population genetics has aspired to build interpretable models of spatial population
31 dynamics (Guillot *et al.*, 2009; Barton, Etheridge and Véber, 2013). These include classic

40 theoretical models such as Fisher's wave-of-advance model (Fisher, 1937), Wright's isolation-by-
41 distance model (Wright, 1943), Kimura's stepping-stone model (Kimura, 1953; Kimura and Weiss,
42 1964), and Malecot's lattice model (Malécot, 1951; Nagylaki, 1976; Rousset, 1997). The field also
43 has a long history of modeling continuous spatial genetic variation (Levene, 1953; Slatkin, 1973;
44 Barton, 1979; Beerli and Felsenstein, 2001; McRae, 2006; Duforet-Frebourg and Blum, 2014;
45 Bradburd, Coop and Ralph, 2018), inferring spatial covariates associated with genetic patterns
46 (Hanks and Hooten, 2013) and detecting spatial barriers to migration (Safner *et al.*, 2011;
47 Petkova, Novembre and Stephens, 2016; Ringbauer *et al.*, 2018; Al-Asadi *et al.*, 2019; Marcus *et al.*,
48 2021). However, these latter efforts are hampered by a lack of good theoretical predictions for
49 continuous, two-dimensional models (Felsenstein, 1975; Barton, Depaulis and Etheridge, 2002),
50 and simulations can provide a valuable tool in the absence of analytical theory.

51
52 The dramatic increase in the number of published whole-genome sequences in the last 20 years
53 (1000 Genomes Project, 2010; Mallick *et al.*, 2016; Palkopoulou *et al.*, 2018; Feuerborn *et al.*,
54 2021), and the advent of ancient genomics (Green *et al.*, 2010; Rasmussen *et al.*, 2010), have
55 revealed previously unknown migration events in the history of several species, such as dogs
56 (Bergström *et al.*, 2020), horses (Librado *et al.*, 2021), elephantids (Meyer *et al.*, 2017), and
57 humans (Lazaridis *et al.*, 2014; Fu *et al.*, 2016). Since migration of populations involves spatial
58 displacement, populations trace their ancestry to different geographic locations (Ralph and Coop,
59 2013; Osmond and Coop, 2021; Wohns *et al.*, 2022). In the context of human history, processes
60 including past migration, gene flow, and population turnovers have been shown to have had a
61 major influence on the present-day distribution of genomic variation (Pickrell and Reich, 2014;
62 Slatkin and Racimo, 2016). Properly anchoring these past demographic events in both time and
63 space has been a focus for new modeling approaches (Racimo *et al.*, 2020; Osmond and Coop,
64 2021; Wohns *et al.*, 2022), and is a question of high interest not only in genetics (Bradburd and
65 Ralph, 2019) but also in ecology (Frachetti *et al.*, 2017; Loog *et al.*, 2017; Crabtree *et al.*, 2021;
66 Delser *et al.*, 2021).

67
68 Despite the key role of geography in population genetics, tools specifically designed for describing
69 and simulating complex spatio-temporal processes are still lacking. Spatial simulations are
70 important not just for rigorous testing and evaluation of existing inference tools and facilitating the
71 development of new inference methods (Liu *et al.*, 2006; Currat and Excoffier, 2011; Delser *et al.*,
72 2021; Osmond and Coop, 2021; Wohns *et al.*, 2022), but also for gaining intuition about the
73 expected behavior of the processes influencing the patterns of genetic variation under various
74 scenarios of spatial population dynamics (Felsenstein, 1975; Slatkin and Excoffier, 2012). While
75 powerful simulation approaches based on coalescent theory have been developed (Hudson,
76 2002; Ewing and Hermisson, 2010; Staab *et al.*, 2015; Kelleher, Etheridge and McVean, 2016),
77 these have little or no notion of spatiality due to fundamental obstacles to incorporating space into
78 the coalescent framework (Barton, Depaulis and Etheridge, 2002; Barton, Etheridge and Véber,
79 2010, 2013), although recent algorithmic advances are promising (Kelleher, Etheridge and
80 Barton, 2014). The first pioneering attempt at simulating spatial population genetic data was the
81 software package SPLATCHE (Currat, Ray and Excoffier, 2004; Currat *et al.*, 2019). However,
82 SPLATCHE's simulation engine is limited to discrete demes based on the stepping-stone model,
83 allows simulation of no more than two populations co-existing at a time, and is not suitable for

84 simulating sequence data at a whole-genomic scale (Currat *et al.*, 2019). The most advanced
85 simulator with spatial capabilities is currently the forward population genetic simulation framework
86 SLiM (Haller and Messer, 2017, 2019). Highly popular in the population genetics community, SLiM
87 contains a vast library of features for simulating individuals in continuous space (as opposed to
88 older approaches based on discrete demes), including spatial interactions between individuals,
89 neighborhood-based mate selection, and customisable offspring dispersal (Haller and Messer,
90 2019). Moreover, the recent implementation of tree-sequence recording in SLiM has opened up
91 the possibility of efficient simulation of massive genome-scale and population-scale datasets
92 (Haller *et al.*, 2019).

93
94 Despite these advances in population genetic simulations, geospatial data analysis remains a
95 complex field with a steep learning curve. Performing even basic manipulations of spatial
96 cartographic objects, handling diverse data formats, and transforming data between different
97 projections and coordinate reference systems (CRS) requires a non-trivial amount of domain-
98 specific knowledge (Lovelace, Nowosad and Muenchow, 2019). Moreover, because the
99 technicalities of geospatial computation are generally not within the scope of population genetic
100 software, available tools do not provide dedicated functionality for building complex and dynamic
101 spatial population models in a straightforward manner. Developing such models and simulating
102 data from them currently requires hundreds of lines of custom code, which is error-prone and
103 hinders reproducibility. Additionally, the lack of specific frameworks for analyzing and visualizing
104 spatially-explicit genomic data further hinders the methodological and empirical progress in spatial
105 population genetics. A flexible and easy-to-use simulation framework specifically designed for
106 developing spatio-temporal population models and analyzing spatial genomic data would expand
107 the horizons of the field, allowing researchers to evaluate the accuracy of novel spatial methods,
108 to test detailed hypotheses about demography and selection, and to answer entirely new kinds of
109 questions about the interactions between organisms across space and time. For instance, many
110 conceptual models and visualizations of past migration events involve depictions of movements
111 of large population ranges across a map as various environmental or cultural conditions change;
112 however, there is currently no easy way to simulate these movements and generate realistic
113 spatio-temporal genomic data.

114
115 To address these issues, we have developed a new programming framework, called *slendr*,
116 designed for simulating and analyzing spatially-explicit genomic data (available at www.slendr.net
117 with extensive documentation and tutorials). The core component of this framework is an R
118 package which leverages real Earth cartographic data (or, alternatively, an abstract user-defined
119 spatial landscape) to programmatically and visually encode spatial population boundaries and
120 their temporal dynamics across time and space, including expansions, migrations, population
121 splits, and gene flow. Because of the challenges involved in testing and validating complex
122 models, *slendr* encourages an interactive workflow in which each component of the model can be
123 inspected and visualized as the model is incrementally constructed in a “bottom-up” fashion.
124 Spatio-temporal models programmed in *slendr* can then be executed using a SLiM back-end
125 script which is bundled with the package and can be controlled by a dedicated R function without
126 leaving the R environment. Additionally, traditional, random-mating, discrete-deme, non-spatial
127 population models can also be simulated, either in forward time using the aforementioned SLiM

128 script or using an alternative coalescent *msprime* (Baumdicker *et al.*, 2022) back-end script which
129 is also bundled with the R package and can provide a more efficient simulation engine for non-
130 spatial models. Both simulation engines of *slendr* save genomic outputs in the form of an efficient
131 tree-sequence data structure (Kelleher *et al.*, 2018), and the *slendr* R package provides a set of
132 functions for loading and processing tree-sequence output files and computing population
133 statistics on them by seamlessly integrating the *tskit* tree-sequence analysis Python module into
134 its R interface. Additional functionality includes conversion of individual trees to a standard R *ape*
135 phylogenetic format (Paradis and Schliep, 2019), and automatic transformation of spatial tree-
136 sequence table data to the standardized *sf* format for geospatial data analysis in R (Pebesma,
137 2018).

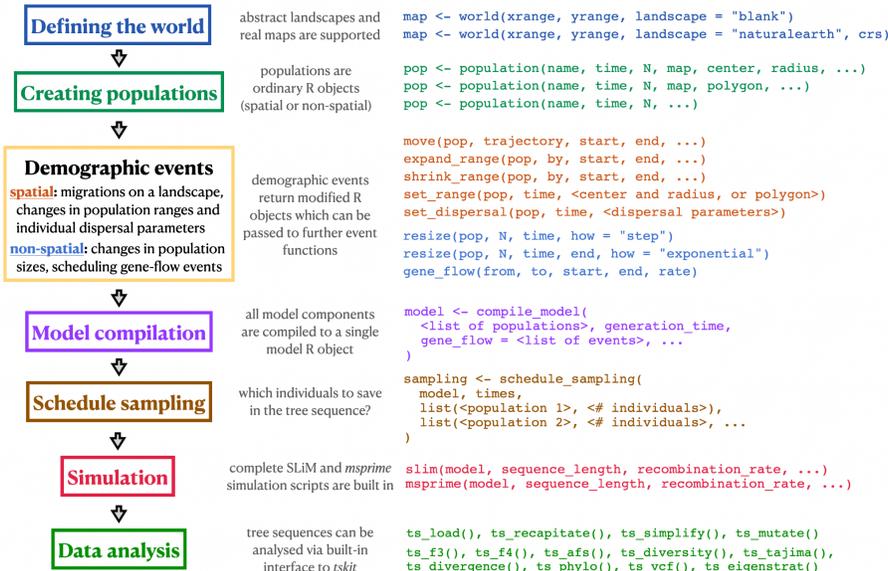
138 Overall, the *slendr* R package facilitates reproducibility by providing a unified framework
139 for writing complete spatial simulation and analysis pipelines entirely in R, which we demonstrate
140 with several concrete examples.

141 Overview of the *slendr* design and typical workflow

142 From a software design perspective, the *slendr* R package represents a tight integration of three
143 distinct parts. First, it implements an interactive and visually-focused R interface for encoding
144 spatio-temporal population dynamics focused on building arbitrarily complex models from small
145 individual components (i.e., simple R objects), designed to require only a minimum amount of
146 code. Second, *slendr* includes two back-end simulation scripts implemented in SLiM (Haller and
147 Messer, 2019) and *msprime* (Baumdicker *et al.*, 2022). These scripts are bundled with *slendr*, are
148 specifically tailored to interpret *slendr* demographic models, and produce tree-sequence files as
149 output (Haller *et al.*, 2019). Lastly, *slendr* provides an interface to the *tskit* tree-sequence analysis
150 library (Kelleher *et al.*, 2018). Although this library is written in C and Python, *slendr* exposes its
151 functionality to the R environment in an R-idiomatic way, blending it naturally with the popular
152 “*tidyverse*” philosophy of data analysis (Wickham *et al.*, 2019).

153 Although these three parts operate at fundamentally different levels under the hood, this
154 integrated approach allows all steps of a *slendr* workflow—from specifying spatio-temporal
155 demographic models, to executing simulations and analyzing simulation results—to be performed
156 without leaving the R environment (**Figure 1**). This allows the user to leverage R’s features for
157 visualization and interactive data analysis at every step of the analytic pipeline, and facilitates
158 reproducibility by eliminating the need to manually integrate disparate software tools and
159 programming languages (Sandve *et al.*, 2013). In this way, *slendr* follows the footsteps of the
160 original design of the S (and later R) languages: to present a consistent and convenient data-
161 analysis-focused domain-specific front end to more efficient and faster tools written in other
162 languages and frameworks (in this case SLiM and *msprime*) (Chambers, 2020).

163
164 In the remainder of this section we outline the individual steps of a typical *slendr* simulation and
165 analysis workflow, as well as describe the individual building blocks of the three main components
166 of the *slendr* framework mentioned above.
167



168
169
170
171
172
173
174
175
176
177
178

Figure 1. Schematic overview of a hypothetical *slendr* simulation and analysis workflow. The colored rectangles on the left indicate individual steps of a hypothetical *slendr* workflow. Short code snippets in matching colors on the right show examples of *slendr*'s declarative interface used in each step, focusing only on a selected few relevant functions and their most important arguments (additional optional arguments are replaced by the "..." ellipsis symbol). The full function reference index can be found at slendr.net/reference. Note that regardless of whether a spatial or non-spatial *slendr* model is being defined and simulated, the workflow remains identical: the same functions are used for both types of models, and the spatial or non-spatial nature of a model is automatically detected by *slendr*.

179 **Defining the world**

180 At the beginning of a *slendr* workflow, the user defines the parameters of the world that the
 181 simulation will occupy using the function `world()` (Figure 1). If the simulated world represents
 182 a region on Earth, the appropriate set of vectorised spatial features will be automatically
 183 downloaded from a public-domain cartographic database (www.naturalearthdata.com). The user
 184 can also specify a dedicated coordinate reference system (CRS) appropriate for the projection of
 185 the geographic region of interest in order to minimize the distortion of distances and shapes
 186 inherent to transforming geometries (in this case population ranges and landscape features) from
 187 the three-dimensional Earth surface to its two-dimensional representation on a map. Alternatively,
 188 the world can be represented by an abstract landscape, optionally with custom features such as
 189 islands, barriers, or corridors. If a non-spatial deme-based model is to be simulated, this step can
 190 be omitted and no changes to the downstream steps described below are needed.

191 **Creating populations and scheduling demographic events**

192 Populations in *slendr* are created with the `population()` function which creates a simple R
193 object containing the parameters of the population that was created (**Figure 1**). In addition to
194 specifying the name, time of appearance, and initial number of individuals for the new population,
195 the user can also specify a `world` object and, if desired, a set of coordinates for the spatial range
196 that the population will occupy. For convenience, the coordinates of all spatial objects in *slendr*
197 (maps, geographic regions, population ranges) are always specified in the global geographical
198 CRS (i.e., degrees of longitude and latitude) but are then automatically internally transformed into
199 the chosen projected CRS (which uses units of meters) if it was specified when creating the
200 `world` (**Figure 1**). This way, users can encode spatial coordinates in familiar units of longitude
201 and latitude while *slendr* internally maintains the proper shapes and distances of spatial features
202 by performing all spatial transformations in the projected CRS.

203 All *slendr* spatial objects are internally represented using a data type implemented by the
204 R package *sf* (Pebesma, 2018), which has emerged as the *de facto* standard for geospatial data
205 analysis in R (Lovelace, Nowosad and Muenchow, 2019). Despite the convenience of the *sf*
206 framework, manipulation of geospatial objects in *sf* still requires writing a non-trivial amount of
207 code dealing with low-level technical details (manipulating and transforming the coordinates of
208 points, lines and polygons). Because most of these technical details are not relevant for specifying
209 population genetic models, we designed a set of domain-specific functions for encoding spatial
210 population dynamics which are expressed in terms of population genetics concepts rather than
211 geometric transformations (**Figure 1**). For instance, the `move()` function accepts a *slendr*
212 population object (i.e., internally an *sf* object, encapsulating the low-level geometric coordinates
213 of the population), a trajectory given as a list of coordinates in longitude and latitude, and a
214 timespan over which the population displacement should occur (*Example 3* and **Figure 4**). Other
215 kinds of dynamic spatial events (population range expansions and contractions, for example) are
216 implemented in an analogous manner. Other demographic events, such as population size
217 changes and gene flow, can be scheduled similarly with another set of straightforward functions
218 (**Figure 1**).

219 For spatial models, the user has the option to fine-tune the within-population individual
220 dispersal and mating dynamics (described in detail in *Example 2* and **Figure 3**) using a set of
221 parameters such as the maximum mating distance between individuals, the dispersal distance of
222 offspring from their parents (and the kernel function of this dispersal), or the parameter influencing
223 the uniformity of the dispersal of individuals within their population's spatial boundary. These can
224 be assigned for each population separately or kept at their default values given in the
225 `compile_model()` step (as we show in *Example 2*). The `competition` parameter determines
226 the maximum neighborhood distance in which individuals in a SLiM simulation compete with each
227 other for space. If this distance is small, then individuals with nearby neighbors have much lower
228 fitness. If the distance is larger, then the effects of crowding are more diffuse. However, if this
229 distance is larger than the dispersal distance (as in *Example 2*), populations tend to self-organize
230 into an evenly-spaced grid of patches. (**Figure 3C**). Using the `competition` parameter, within-
231 population dynamics can thus be fine-tuned to represent various levels of individual clustering
232 into sub-groups (**Figure 3C**). In addition to the `competition` parameter, a `mating` parameter
233 determines the maximum distance to which an individual will look for a mate to produce offspring.

234 Finally, a `dispersal` parameter determines how far an offspring can end up from its parent, and
235 a related `dispersal_fun` argument characterizes the density function for this dispersal:
236 "normal" (default), "uniform", "cauchy", "exponential", or "brownian"; more details are available in
237 the *slendr* R package documentation at slendr.net/reference. We note that changes in all three
238 spatial interaction and dispersal parameters can be also scheduled dynamically at specific times
239 throughout the run of a model with a *slendr* function `set_dispersal()`.

240
241 A standard feature of many population genetic frameworks is the specification of the times of
242 various demographic events in terms of generations, either forwards in time starting from
243 generation 1 (as is the case with SLiM) or backwards in time starting from time 0 "in the present"
244 (as is the case with coalescent frameworks such as *msprime*). This can be cumbersome in cases
245 when the events or samples of interest are traditionally specified in times of "years before present"
246 (such as dated ancient DNA samples), or in situations in which it would be desirable to simulate
247 future outcomes, as in ecological predictive modeling. Moreover, because these standard times
248 often need to be converted into generations by a factor specifying the length of the generation
249 time of the species of interest, this can easily lead to frustrating bugs in simulation scripts. To
250 ameliorate this situation, *slendr* allows the users to specify times in whichever time units they
251 would prefer, in either the forward or backward direction. The time direction is automatically
252 detected by *slendr* from the sequence of demographic events specified for a model (but can also
253 be set explicitly), and the conversion of event times into generations is performed in the
254 compilation step via the provided `generation_time` argument to `compile_model()`
255 (described below). Similarly, times of the tree-sequence nodes in *slendr*'s outputs (which are
256 specified by most simulation software in terms of generations backwards in time) are
257 automatically converted by *slendr* back into the units of time used by the user during model
258 specification.

259
260 Because every *slendr* demographic event function returns a modified population object which can
261 be further used as an input to other *slendr* functions, the R interface encourages a workflow in
262 which complex models are composed incrementally from smaller components (**Figure 1**,
263 *Examples 1-3*). Importantly, because each *slendr* function assures the consistency of the model
264 by enforcing appropriate constraints during the model definition process (e.g., a population cannot
265 be moved or participate in a gene-flow event at a time when it would not yet exist), this workflow
266 facilitates the early discovery of bugs before the simulation (which can be extremely
267 computationally costly) is even executed. This is further facilitated by a convenient set of plotting
268 functions, such as `plot_map()` and `plot_model()`, which can visualize the spatio-temporal
269 dynamics of the specified model (or its individual components) as the model is being incrementally
270 developed.

271 Model compilation

272 Having defined all the individual components of a population model (i.e., created all the necessary
273 population and gene-flow events), the user calls the function `compile_model()` to compile the
274 model configuration to a single R object (**Figure 1**)—a step in which *slendr* performs additional
275 checks for model consistency and correctness. Furthermore, this operation also transforms the

276 model components from their R representation into a set of files on disk, written in a format
277 interpretable by the built-in SLiM and *msprime* simulation back-end scripts which are used to
278 execute *slendr* models in the next phase, as described below. The compiled model object can
279 also be used as input for a built-in R-based interactive browser app built using the *shiny* R
280 package (Chang *et al.*, 2021) which allows the user to “play” the defined spatial model dynamics
281 over time and explore the “admixture graph” implied by the model (Patterson *et al.*, 2012) for
282 additional verification of the model’s correctness. The functions `plot_map()` and
283 `plot_model()` mentioned above also accept a compiled model object as their input and produce
284 a static visualization of the model.

285 Scheduling sampling events and simulation

286 The *slendr* package comes bundled with two simulation back-end scripts which were tailored to
287 interpret the configuration files produced by the `compile_model()` function and simulate the
288 model, triggering all of the encoded population dynamics in the course of the simulation run.

289 The first back-end script is written in SLiM’s programming language Eidos (Haller and
290 Messer, 2019), and can execute both spatial and non-spatial *slendr* models in a Wright–Fisher
291 setting by calling *slendr*’s `slim()` function. The second back-end script is implemented using
292 *msprime* (Baumdicker *et al.*, 2022) and is designed to interpret the compiled *slendr* model in a
293 non-spatial setting as a standard coalescent simulation by calling *slendr*’s `msprime()` function.
294 Both simulation engines can interpret the same *slendr* model without a need to make any
295 changes. For instance, a spatial model can be run with the *msprime* back end, in which case the
296 spatial component of the model is simply ignored. Because coalescent simulations are generally
297 much more computationally efficient than their forward-time counterparts, the *msprime* back end
298 of *slendr* can be useful for R users who would like to run a large number of traditional, non-spatial
299 simulation replicates efficiently without having to write custom Python *msprime* code or use its
300 *ms*-like command-line interface (Hudson, 2002). Importantly, the correctness of both *slendr*
301 simulation engines is validated using a set of automatic statistical tests on non-spatial models
302 which ensure that when a *slendr* model is run in both SLiM and *msprime*, the demographic events
303 specified by the model (population splits, population size changes, and gene-flow events) result
304 in equivalent site-frequency spectra and *f*-statistics (Patterson *et al.*, 2012) between both back
305 ends.

306 Leveraging the ability to save simulation outputs as a tree sequence (Kelleher *et al.*, 2019;
307 Speidel *et al.*, 2019) from both SLiM (Haller *et al.*, 2019) and *msprime* (Baumdicker *et al.*, 2022),
308 *slendr* embraces the tree sequence as its primary output format. This is powerful not only because
309 the tree sequence represents an extremely efficient representation of even large-scale population
310 genomic data, but also because it provides an elegant way to calculate many population genetic
311 statistics of interest, a feature which we describe in more detail in the next section. To specify
312 which simulated individuals should be recorded in the output tree sequence, *slendr* provides two
313 alternative approaches. First, if no explicit sampling schedule is specified, all individuals living at
314 the very end of a SLiM simulation run are explicitly sampled (i.e., “remembered”) in the tree
315 sequence output, matching the default behavior of SLiM. If a *slendr* model is simulated with the
316 *msprime* back end, the number of recorded individuals will be equal to the population size of each
317 population at the start of the coalescent process looking backwards in time (i.e., in “the present”).

318 Alternatively, *slendr* provides a flexible way to trigger sampling events via its
319 `schedule_sampling()` function, which allows one to specify the time (and, optionally, the
320 location) at which a sample comprising a given number of individuals from a given population
321 should be taken and recorded in the tree sequence (*Example 3*). To improve readability and
322 interpretation of *slendr* analysis code, every sampled individual can be referred to using its
323 readable name during tree-sequence processing and computation of statistics (*Examples 1, 2,*
324 *and 4*) rather than just by numeric identifiers as is the case with the default tree-sequence analysis
325 workflow with *tskit* (Kelleher *et al.*, 2018).

326 Data analysis

327 The default output of a *slendr* simulation is a tree sequence. However, because processing and
328 analysis of tree-sequence files requires a non-trivial knowledge of Python or C (Kelleher *et al.*,
329 2018) which many R users might not have, *slendr* provides an R-idiomatic interface to the most
330 commonly used *tskit* tree-sequence methods such as the allele frequency spectrum, Patterson's
331 *f*-statistics, and various summary statistics of population diversity (Patterson *et al.*, 2012; Ralph,
332 Thornton and Kelleher, 2020). This way, users can design population genetic models in R,
333 execute them from R using the built-in `slim()` or `msprime()` functions, and analyze the
334 resulting tree sequence data without having to leave the R environment for downstream statistical
335 analyses and plotting, and without the need to convert outputs to other bioinformatic or population
336 genetic file formats. Although primarily designed for analysis of tree sequences generated from
337 *slendr* models, the R-*tskit* interface can operate also on tree sequences without *slendr*-specific
338 metadata. Therefore, users who would prefer to run simulations with standard *msprime* or SLiM
339 scripts but are interested in analyzing their tree-sequence results in R will still find the *slendr* R
340 package useful. The reference manual at slendr.net/reference contains a complete list of *tskit*
341 tree-sequence methods that have been integrated into *slendr*'s R interface. If integration with
342 traditional tools such as PLINK (Purcell *et al.*, 2007) or ADMIXTOOLS (Patterson *et al.*, 2012) is
343 required, functions for exporting to VCF (Danecek *et al.*, 2011) and EIGENSTRAT (Patterson *et*
344 *al.*, 2012) are also provided.

345 During a spatial simulation in SLiM, each sampled individual's location on the simulated
346 landscape is tracked and recorded in the tree sequence, encapsulating the full spatio-temporal
347 genealogical history that has been simulated. When the tree-sequence output file is then loaded
348 by *slendr*, *slendr* processes the spatial locations of nodes in the tree sequence (which represent
349 chromosomes of past and present individuals), and transforms them back into the original
350 coordinate system of the simulated world, adding additional annotation data such as readable
351 names of sampled individuals, population assignments of each individual and node, etc.
352 Furthermore, this information is exposed in an *sf*-compatible format, meaning that the spatio-
353 temporal information about ancestral relationships between simulated samples can be processed,
354 analyzed, and visualized using a wide range of R packages including *sf*, *ggplot2*, and *dplyr*
355 (Pebesma, 2018; Wickham *et al.*, 2019). Additionally, individual trees in the tree sequence can
356 be extracted by a *slendr* function, `ts_phylo()`, which converts *tskit*-formatted tree objects into
357 the format defined by the R phylogenetics package *ape*, which has been the standard for
358 phylogenetics in the R ecosystem for nearly two decades (Paradis and Schliep, 2019). This gives

359 *slendr* users even more options to analyze tree-sequence results with a large array of standard
360 phylogenetics tools available for the R environment (Paradis, 2011).

361 Installation and software dependencies

362 *slendr* is currently developed for macOS and Linux. It is available on the CRAN R package
363 repository at <https://CRAN.R-project.org/package=slendr>, and can be installed from the
364 interactive R console with the standard command `install.packages("slendr")`.
365 Development versions of *slendr* which contain latest bug fixes and new experimental features can
366 be installed from its GitHub repository using the R package *devtools* with the R command
367 `devtools::install_github("bodkan/slendr")`.

368 Two external software dependencies must be present on a user's system to leverage the
369 full functionality of *slendr*: a forward population genetic simulator SLiM (Haller and Messer, 2019)
370 (which is required for running spatial simulations and non-spatial simulations in the forward-time
371 setting) and a trio of Python modules *msprime* (Baumdicker *et al.*, 2022), *tskit* (Kelleher *et al.*,
372 2018) and *pyslim* (github.com/tskit-dev/pyslim) (which are needed to run *slendr* models as
373 coalescent simulations and to analyze tree-sequence data).

374 The SLiM software is available for all major operating systems and its installation
375 instructions can be found at messerlab.org/slim. Importantly, the current version of *slendr* requires
376 the latest release of SLiM 4.0. In order to use SLiM for simulations in *slendr*, the R session needs
377 to be aware of the path to the directory containing the SLiM binary. Calling `library(slendr)`
378 for the first time provides an informative message for the user on how this can be accomplished
379 by modifying the `$PATH` variable by editing the `~/.Renviron` file.

380 Because some users might find the experience of setting up a dedicated Python
381 environment with the necessary Python modules challenging (especially users who exclusively
382 work with R), *slendr* provides an R function `setup_env()` which automatically downloads a
383 completely separate Python distribution and installs the required versions of *tskit*, *msprime*, and
384 *pyslim* Python modules in their correct required versions into a dedicated virtual environment
385 without any need for user intervention. Moreover, this Python installation and virtual environment
386 are isolated from other Python configurations that might be already present on the user's system,
387 thus avoiding potential conflicts with the versions of Python and Python modules required by
388 *slendr*. Once this isolated Python environment is created by `setup_env()`, users can activate it
389 in future R sessions by calling a helper function `init_env()` after loading *slendr* via
390 `library(slendr)`. Therefore, although *slendr* uses Python modules for internal handling of
391 tree-sequence data and coalescent simulation, direct interaction with Python is not necessary.
392

Deleted: After a dedicated Python environment is created by `setup_env()`, calling `library(slendr)` at any later point will activate this environment automatically. ...

393 [Relationship of *slendr* to SLiM and *msprime*](#)

394 [Given that *slendr*'s simulation engines are implemented in SLiM and *msprime*, it is worth
395 elaborating on its relationship to these simulation frameworks, particularly in terms of the features
396 supported by *slendr*. First, it is important to note that *slendr* is not simply a wrapper for SLiM and
397 *msprime* in the strict sense of the word, since *slendr* does not provide an R equivalent of every
398 function and method provided by SLiM and *msprime*. Instead, *slendr* aims to provide a user-](#)

403 [friendly, R-idiomatic way to encode a particular class of “traditional” Wright-Fisher population](#)
404 [genetic models frequently used in evolutionary biology and population genetics, allowing users to](#)
405 [employ such models with a minimal amount of coding. Most importantly, *slendr* models currently](#)
406 [assume that populations evolve via random mating, and that the genomes of individuals evolve](#)
407 [neutrally, with mutations overlaid on top of the simulated genealogies after each simulation run.](#)
408 [This applies also to spatial *slendr* demographic models, with the caveat that interaction and](#)
409 [dispersal distance parameters can—depending on the exact parametrization of each spatial](#)
410 [slendr model—cause individuals to only mate locally, which can have interesting implications for](#)
411 [the behavior of standard population genetic statistics \(as shown in *Example 2*\).](#)

412 [The complete set of models supported by *slendr* is likely to slightly expand over time as](#)
413 [new features are implemented. Details of new features, such as customized recombination maps](#)
414 [and non-neutral mutation types, are being discussed with the community on the GitHub page of](#)
415 [slendr \(<https://github.com/bodkan/slendr>\), and users are encouraged to provide feedback there.](#)
416 [The four practical examples \(*Examples 1–4* below\) have been designed to demonstrate the full](#)
417 [range of *slendr*’s features at the time of writing.](#)

418 [Finally, because *slendr*’s forward and coalescent simulation back ends are implemented](#)
419 [as fairly standard SLiM and *msprime* scripts, the performance of *slendr* simulations and tree-](#)
420 [sequence analyses can be assessed using already-existing benchmarks and guidelines provided](#)
421 [by publications describing SLiM and *msprime* \(Haller *et al.*, 2019; Baumdicker *et al.*, 2022; Haller](#)
422 [and Messer, 2022\).](#)

423 Practical examples

424 In the following sections, we present the features of the *slendr* R package with several practical
425 examples, each of which focuses on a different aspect of the *slendr* simulation framework. We
426 start by showing how traditional, non-spatial, random-mating models can be specified with a
427 minimum amount of R code (*Example 1*). We then proceed with two examples of spatial models:
428 first, a model showing how the degree of the spatial spread of a population can be adjusted by
429 setting the within-population individual-based dispersal dynamics (*Example 2*); second, a model
430 which schedules the movements of entire population ranges across a landscape (*Example 3*).
431 These examples are intended to demonstrate *slendr*’s ability to define complex spatio-temporal
432 models incrementally, building them from simpler components. We also emphasize how *slendr*
433 model configuration and simulation steps naturally flow into data analysis, all within the R
434 environment. In the final demonstration (*Example 4*), we tap into the rich information embedded
435 in spatial tree sequences to visualize individual trees on a landscape, tracing the complex spatio-
436 temporal ancestry of an individual on the simulated map. Extended versions of these and many
437 other examples with complete reproducible code for simulation, analysis, and plotting can be
438 found as standard R package vignettes at *slendr*’s website (www.slendr.net).
439

440 Example 1: Traditional non-spatial model

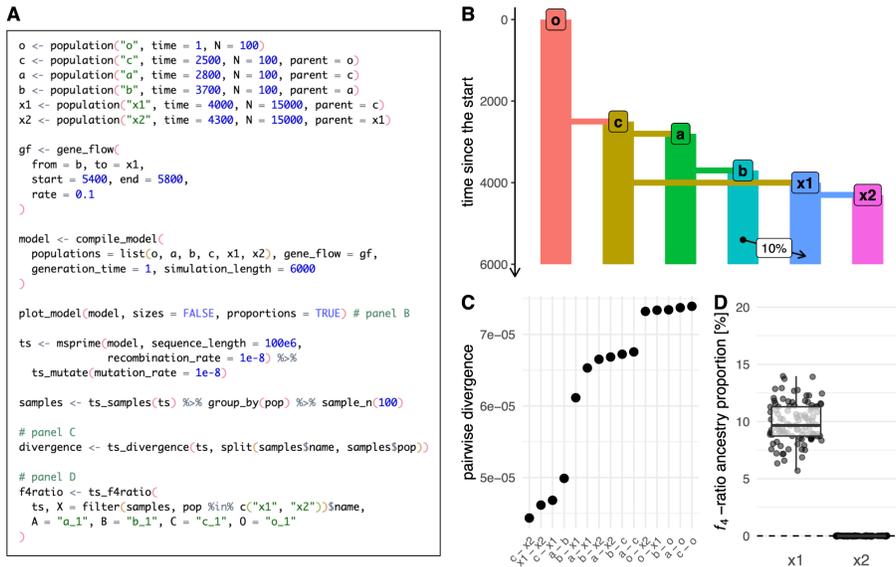
441 Regardless of whether a spatial or non-spatial model is defined and simulated, the *slendr*
442 workflow remains the same. Therefore, before we explore spatial models, we begin by showing

443 how a traditional, non-spatial population genetic model can be constructed with *slendr* and how
444 users can compute population genetic statistics on simulated tree-sequence outputs using
445 *slendr*'s R interface to the *tskit* tree-sequence analysis library (Kelleher *et al.*, 2018) (represented
446 by functions with the `ts_*` prefix, **Figure 1**).

447 First, we define an abstract demographic model similar to that which is commonly used in
448 teaching the principles behind the f_4 -ratio ancestry proportion estimator (Patterson *et al.*, 2012).
449 In *slendr*, we define the model with a straightforward sequence of `population()` calls that
450 schedule the order of splits for several populations, taking care of parent–daughter population
451 relationships by providing the appropriate population object as a `parent` argument when creating
452 each daughter population (**Figure 2A**). We then schedule a single gene-flow event between the
453 populations “*b*” and “*x1*” by calling the `gene_flow()` function. After compiling the model with
454 `compile_model()`, we verify its correctness by visualizing the embedded population
455 relationships with `plot_model()` (**Figure 2B**). Although only a single `gene_flow()` event is
456 featured in this example, more complex gene-flow networks can be specified with *slendr*.
457 Conveniently, strict consistency checks validate each encoded gene-flow event before the
458 computationally costly simulation is run. Examples of complex models with dozens or hundreds
459 of gene-flow events can be found in the documentation available on the *slendr* website
460 (www.slendr.net).

461 As stated before, *slendr* provides two simulation back ends; here we use the coalescent
462 *msprime* back end to simulate the model, since SLiM's spatial capabilities are not required for this
463 simple non-spatial model. However, we note that the function `slim()` could be used in place of
464 the `msprime()` call to perform the equivalent forward-time simulation just as easily. By default,
465 *slendr* automatically loads the simulated tree-sequence object which can be immediately used for
466 analysis. In this example, we compute the pairwise divergence between random samples of 100
467 individuals from each population with the function `ts_divergence()` (**Figure 2C**). Finally, we
468 use the function `ts_f4ratio()` to compute the values of the f_4 -ratio estimate of “*b*” ancestry in
469 populations “*x1*” and “*x2*”, which differ in whether or not they experienced gene flow from “*b*”
470 (**Figure 2D**). All other tree sequence analysis functions of *slendr* (**Figure 1**) can be accessed in
471 the same way. We note that because *slendr* assigns symbolic, permanent names to individuals
472 during sampling, the users can refer to them with these names during tree-sequence operations
473 such as simplification and when computing tree-sequence statistics.

474
475



476
477
478
479
480
481
482
483
484
485
486
487
488
489
490

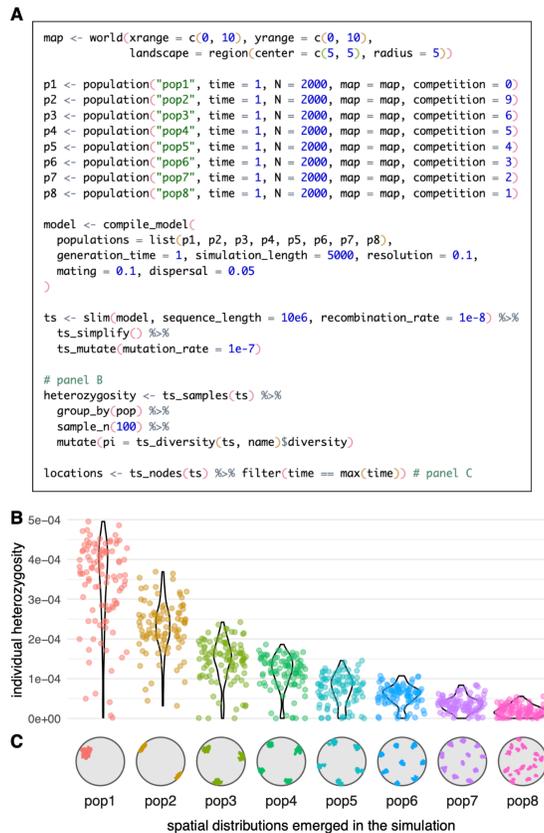
Figure 2. Example 1: specifying a non-spatial model and computing statistics on tree-sequence output. (A) A script which defines a model of a simple demographic history of six populations, simulates it with the *msprime* back end by calling the function `msprime()`, and performs analyses shown in B–D. (B) A visual overview of the compiled *slendr* model produced by `plot_model()` prior to simulation. (C) Visualization of the data frame produced by `ts_divergence()` on the output tree sequence simulated. (D) Ancestry proportions estimated with `ts_f4ratio()` directly from the tree sequence output. As expected from the model definition, the f_4 -ratio statistic estimates indicate ~10% ancestry from “b” in the population “x1”, but 0% ancestry in population “x2”; this agrees with the model overview shown in panel B. Full *ggplot2* visualization code for the figures can be found in a vignette dedicated to this paper at www.slendr.net. The runtime for the simulation and analysis shown in A was ~5 minutes, as measured on a 16” MacBook Pro (2021) equipped with the Apple M1 Pro chip, 32 GB RAM, and running macOS Ventura 13.1.

491 **Example 2: Model with population dispersal dynamics**

492
493
494
495
496
497
498
499

In our second example, we move from a non-spatial, random-mating model to a model which is explicitly spatial. First, we create an abstract, circular world map using the function `world()`, producing a completely featureless landscape (see *Example 3* for a more elaborate world map). We then create a series of eight populations which all occupy that map, as specified by the `map` argument to `population()`, but do not interact with each other. For simplicity, each population forms its own evolutionary lineage without additional splits or gene-flow events. Importantly, we set the `competition` parameter of each population to a value which forces the individuals to

500 assume an increasing degree of spatial subdivision which, in turn, affects the amount of diversity
501 expected in each population. Finally, we compile the model to a single object with
502 `compile_model()` and run it with the `slim()` back end, simulating 16,000 diploid genomes of
503 10 megabases each (**Figure 3A**). After the simulation finishes, we simplify the produced tree
504 sequence, overlay mutations on the simulated genealogies, and use the `slendr` function
505 `ts_diversity()` to compute the expected heterozygosity in a sample of 100 individuals from
506 each population, inspecting how heterozygosity is affected by the emergent spatial arrangement
507 of each population (**Figure 3B, C**). We note that some of the values of the spatial competition
508 distance parameter used in this example are quite large, especially compared to the much shorter
509 maximum distance of individual dispersal and mating. Although biologically rather unrealistic, the
510 competition distances have been chosen to give rise to very different degrees of spatial
511 subdivision and, consequently, to varying levels of population genetic diversity, with the intention
512 to demonstrate the ease with which a wide range of model dynamics can be configured by the
513 user.
514



515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530

Figure 3. Example 2: a spatial model which involves the parametrization of within-population dispersal dynamics. (A) A complete script which defines eight populations as independent lineages or species, each with constant size and each defined with a different value of *slendr*'s spatial competition parameter, with analysis code to produce panels **B–C**. The simulation is run with *slendr*'s SLiM back end for 5000 generations, after which a tree sequence recording the genealogical history of 2000×8 diploid individuals is loaded, simplified, and mutated. Heterozygosity is then computed for 100 individuals randomly sampled from each population at the end of simulation. **(B)** Distribution of heterozygosities of individuals observed in all eight populations. **(C)** A snapshot of the spatial distributions which emerged as a result of the competition parameter value set for each population. Full visualization code for the figures can be found in a vignette dedicated to this paper at www.slendr.net. [The runtime for the simulation and analysis shown in A was ~12 minutes, as measured on a 16" MacBook Pro \(2021\) equipped with the Apple M1 Pro chip, 32 GB RAM, and running macOS Ventura 13.1.](#)

531 Example 3: A toy model of movements and expansions of human
532 populations in West Eurasia over the last 50,000 years

533

534 In this example we further expand on the *slendr* functionality demonstrated in the first two
535 examples, introducing programming of expansions and migrations of entire population ranges
536 across a realistic landscape—perhaps the most distinctive feature of *slendr*. The model we
537 implemented here is inspired by large-scale population migrations and turnover events inferred
538 from ancient DNA analyses of human remains from across West Eurasia (Lazaridis *et al.*, 2014;
539 Allentoft *et al.*, 2015; Haak *et al.*, 2015), although we caution that it is simplified and intended only
540 as an illustrative example.

541 Similarly to *Example 2*, we begin by defining a world map for the simulation (**Figure 4A**),
542 in this case using realistic Earth cartographic data provided by the Natural Earth project
543 (naturalearthdata.com). Because we focus on the broad region of West Eurasia, we select the
544 most appropriate coordinate reference system (CRS) for projecting this region on a two
545 dimensional map which is EPSG:3035. We then define a series of populations, specifying their
546 approximate geographic ranges using simple polygons. We then use the functions `move()` and
547 `expand_range()` to schedule when and where populations should migrate, and by what
548 distance and how quickly their population ranges should expand across the landscape during
549 simulation. We again use `plot_model()` to visualize the demographic history embedded in the
550 *slendr* model as a non-spatial tree-like structure with gene-flow edges (**Figure 4B**); here, we also
551 use `plot_map()` to get a “compressed” overview of the spatio-temporal population range
552 dynamics on the simulated map (**Figure 4C**). We note that unlike in the two previous examples,
553 which were specified in forward time units, this example expresses the timing of demographic
554 events in units of “years before present” which is more natural to this model.

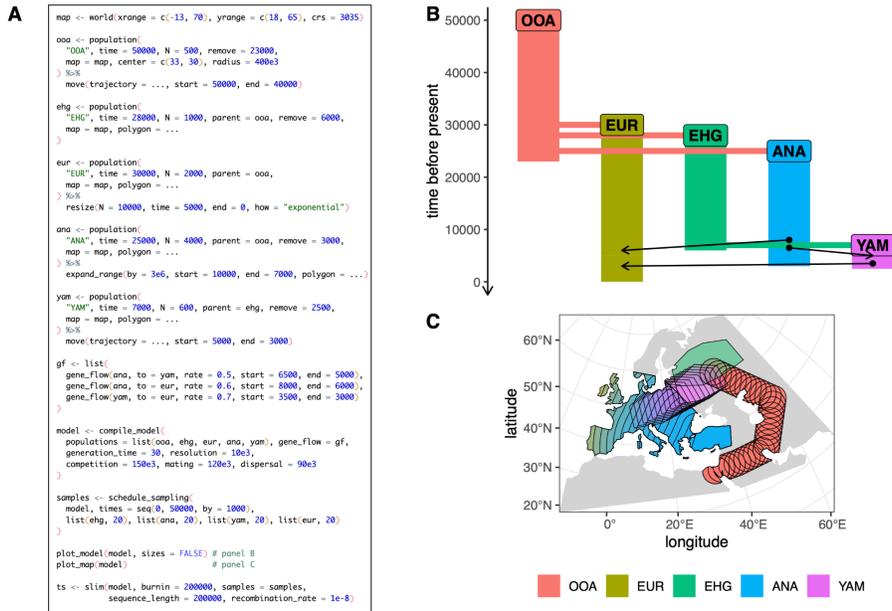
555 In the previous two code examples (**Figure 2A, 3A**) we used the default tree-sequence
556 sampling of *slendr*, which implicitly records the genomes of all the diploid individuals alive at the
557 end of a simulation. In this example, we instead use `schedule_sampling()` to specify a series
558 of sampling events from each population every 1,000 years. We then execute the compiled model
559 and the sampling schedule specified using the `slim()` back-end, which records only the
560 scheduled set of sampled individuals in the tree-sequence output file.

561

562

563

564



565
566
567
568
569
570
571
572
573
574
575
576
577
578

Figure 4. Example 3: a demographic model on a real Earth landscape. (A) A *slendr* script which defines a toy spatio-temporal model of human prehistory in West Eurasia, with analysis code that produces panels **B–C**. For brevity, we do not specify the full set of coordinates for each spatial demographic event or population range polygon, instead indicating them as "..."; the complete reproducible code can be found in a vignette dedicated to this paper at www.slendr.net. **(B)** Visual summary of the non-spatial component of the demographic model, produced by `plot_model()` with arrows indicating gene flow events. **(C)** A “compressed” view of spatio-temporal snapshots of population ranges throughout the course of the model prior to the simulation, produced by `plot_map()`. [The runtime for the simulation shown in A was ~3 minutes, as measured on a 16” MacBook Pro \(2021\) equipped with the Apple M1 Pro chip, 32 GB RAM, and running macOS Ventura 13.1.](#)

579 Example 4: Visualization of individual trees and spatio-temporal ancestral
580 lineages across a landscape

581

582 In our final example (**Figure 5**), we return to the abstract toy model of West Eurasian prehistory
583 developed in *Example 3*. To leverage *slendr*’s power to simulate genomic data from complex
584 spatial demographic models, *slendr* makes it easy to tap into the large library of geospatial data
585 science packages available for R (Lovell, Nowosad and Muenchow, 2019) by automatically

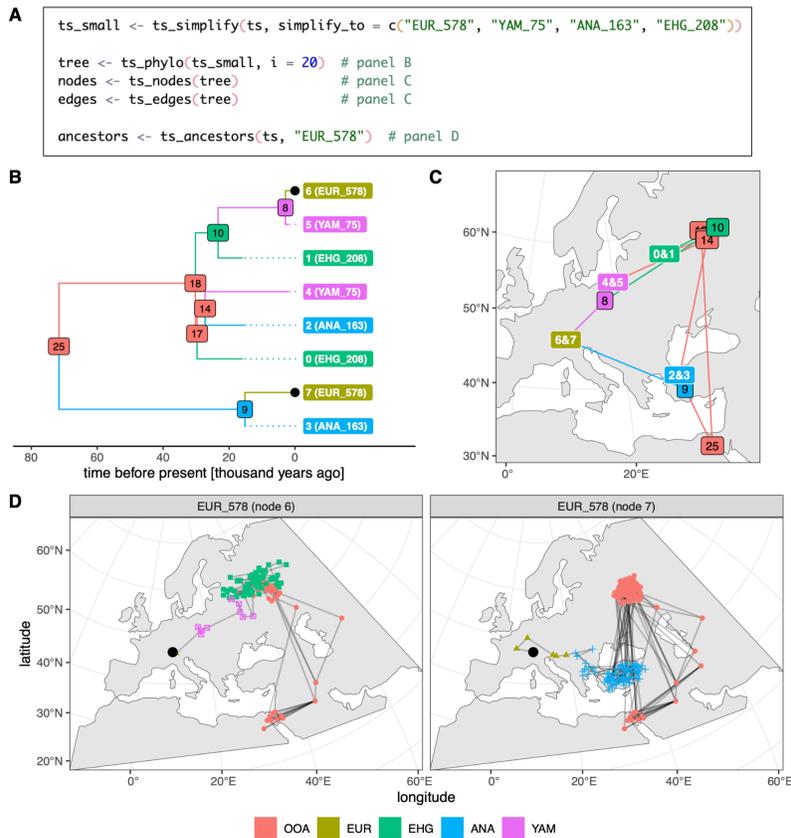
586 converting simulated spatial locations to an *sf*-compatible tabular format (Pebesma, 2018), as we
587 will see here.

588 To demonstrate the richness of the spatio-temporal information recorded in the tree
589 sequence, we use the full tree sequence produced by the code in **Figure 4A** and simplify it so
590 that it contains only the history of a small subset of the thousands of individuals sampled during
591 the spatio-temporal simulation (**Figure 5A**). We then extract the 20th tree in the tree sequence
592 with *slendr*'s function `ts_phylo()`, which converts a tree from the *tskit* tree sequence into an R
593 `phylo` format defined by the *ape* R package, a standard tool for phylogenetics in R (Paradis and
594 Schliep, 2019). Such tree objects can be analyzed by any of the dozens of R packages which
595 operate on *ape*'s `phylo` trees—for instance, in **Figure 5B** we show a visualization of this tree
596 using the R package *ggtree* (Yu *et al.*, 2017). Furthermore, because the tree was generated from
597 a spatially-annotated tree sequence, the user can extract information about the location of each
598 individual (or node) in the tree across space and time, as well as ancestral relationships between
599 nodes in the tree, using `ts_nodes()` and `ts_edges()` respectively. Crucially, because these
600 functions automatically convert locations into *sf*'s geospatial representation (including the
601 appropriate CRS projection), the results can be immediately plotted on a map with *ggplot2*, which
602 has built-in support for *sf* data (**Figure 5C**).

603 In addition to extracting and visualizing single trees representing a genealogy of a set of
604 sampled genomes descending from a common ancestor (spatial or non-spatial), *slendr* also
605 provides a way to extract the complete spatio-temporal ancestry of a single sample going back in
606 time across the entire tree sequence, potentially spanning many trees with thousands of the
607 sample's ancestors. This can be accomplished with the function `ts_ancestors()` which, in an
608 analogous way to `ts_nodes()` and `ts_edges()`, exposes the spatio-temporal information in
609 the tree sequence as an *sf* object which can be visualized on a map with *ggplot2*. In this example,
610 we use `ts_ancestors()` to reconstruct the spatio-temporal ancestry distribution for a single
611 simulated European individual ("EUR_578", represented by the black dot in **Figure 5D**). Because
612 this individual is diploid, we can trace the ancestry carried by its one chromosome through an
613 expansion from Anatolia (**Figure 5D**, right panel), while its other chromosome clearly traces its
614 ancestry to a population which migrated to central Europe from an eastern population (**Figure**
615 **5D**, left panel).

616 Note that by default, the tree sequence output of a *slendr* simulation only contains
617 information about ancestors which are represented by coalescent nodes in some marginal tree—
618 i.e., nodes which are a most recent common ancestor of some pair of sampled nodes. In this
619 example, in **Figure 5B** and **C** we can see that the most immediate ancestor (node number 9) of
620 one chromosome of the sampled individual "EUR_578" lived in the region of Anatolia, but the
621 ancestor of its second chromosome lived in Europe (node number 8); but we do not know where
622 all the ancestors along the edges between nodes 9-7 and 8-6, since they were simplified away.
623 Similarly, **Figure 5D** shows the distributions of locations of most recent common ancestors, not
624 all ancestors. The distribution of ancestors at a particular point in time could be obtained by adding
625 an appropriate sampling event to *slendr*'s sampling schedule and then extracting ancestors from
626 that time.

627



628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643

Figure 5. Example 4: accessing and visualizing spatio-temporal information encoded in trees and tree sequences simulated with the *slendr*. (A) A continuation of the script from Example 3, showing how a (potentially very large) tree sequence generated from a *slendr* model can be simplified to a subset of individuals with `ts_simplify()`. A single tree from the tree sequence is then extracted with `ts_phylo()`, the tables of spatio-temporal locations of nodes and branches of the tree are extracted by `ts_nodes()` and `ts_edges()`, and ancestry information for one individual across the entire tree sequence is extracted with `ts_ancestors()`, in order to produce the data plotted in B–D. (B) A visualization of the tree extracted by `ts_phylo()` using standard visualization features of the *ggplot2* and *ggtree* R packages. Dotted lines indicate shortened branches of ancient samples. (C) Visualization of the tree from panel B as a network across the original spatial simulation landscape, with each node indicating the location of a particular individual who lived at some point during the simulation. Labels with numbers correspond to the locations of sampled individuals, each carrying two chromosomes which are represented by two nodes in the tree sequence. All node numbers correspond to those shown in

644 the tree in panel **B**. The plot was generated with *ggplot2* using the *sf*-formatted data extracted by
645 `ts_nodes()` and `ts_edges()`. **(D)** A visualization of the spatio-temporal ancestry of a single
646 simulated European individual, "EUR_578", using the information from the entire tree sequence.
647 Each sub-panel shows the spatial ancestry distribution of one of the two chromosomes carried by
648 this individual (the location of whom is indicated by a black dot), tracing its ancestry through different
649 lineages all the way back to a population in Africa. For easier reference, the same black dots
650 indicate the two chromosomes of this individual also in the tree in panel **B**. The *ggplot2* code for
651 the figures is omitted for brevity. Full reproducible code examples including the visualization code
652 can be found in a vignette dedicated to this paper at www.slendr.net. [The runtime for the code](#)
653 [shown in A was ~1 second, as measured on a 16" MacBook Pro \(2021\) equipped with the Apple](#)
654 [M1 Pro chip and 32 GB RAM, running macOS Ventura Version 13.1.](#)

655 Discussion

656 The *slendr* R package provides a new programmable framework for simulating complex spatio-
657 temporal genomic data. The package implements a set of features for defining spatial population
658 range dynamics with a declarative and visually-focused R interface and uses a tailor-made SLiM
659 script as an efficient population genetic simulation engine. Additionally, *slendr* provides a
660 convenient new way to simulate and analyze large-scale genomic data sets even from traditional,
661 non-spatial demographic models using *msprime* entirely within the R environment.

662 Owing to its declarative interface, which requires little code even for complex models, the
663 *slendr* package is highly accessible even to researchers or students with little or no prior
664 experience in programming. One of the major challenges for novice population geneticists is
665 having to learn how to integrate multiple different software tools and programming frameworks. R
666 (R Core Team, 2021) is often the first language that biology and bioinformatics students learn,
667 since it offers a large number of libraries for data analysis, statistics, and plotting (Wickham and
668 Grolemund, 2016). For these users, *slendr* provides the opportunity to explore population genetic
669 concepts and simulate realistic population genomic data as soon as they learn the most basic
670 principles of R (i.e., how to call R functions and work with data frames), without first having to
671 learn Python for *msprime* simulations (Baumdicker *et al.*, 2022), shell scripting for simulators from
672 the *ms* family (Hudson, 2002; Staab *et al.*, 2015), or Eidos for SLiM (Haller and Messer, 2019).

673 Tree sequences provide an efficient way to compute many commonly used population
674 genetic statistics directly on the simulated genealogies (Ralph, Thornton and Kelleher, 2020);
675 because *slendr* uses the tree sequence as its default output format (Kelleher *et al.*, 2018; Haller
676 *et al.*, 2019), in many cases users do not need to convert simulation outputs to external file formats
677 such as VCF or EIGENSTRAT for analysis in other software. This way, *slendr* simulations can be
678 readily used in model fitting and population genetic analyses in situations which have traditionally
679 required converting simulated data to genotype files before analyzing them with population
680 genetics tools such as PLINK (Purcell *et al.*, 2007) or ADMIXTOOLS (Patterson *et al.*, 2012). That
681 said, export to VCF and EIGENSTRAT genotype file formats is supported with a single function
682 call (`ts_vcf()` and `ts_eigenstrat()`) if needed.

683 A key principle in the design of *slendr* has been reproducibility (Sandve *et al.*, 2013): a
684 complete *slendr* simulation and analysis workflow can be written as a single R script. Additionally,
685 the compilation of any *slendr* module produces a self-contained "bundle directory" containing all
686 model configuration files and simulation back-end scripts required to execute the model from the

687 command-line. Although accessing this directory is not necessary for standard workflows because
688 *slendr* operates entirely from R, these bundles can be checked into a git history and provided as
689 supplementary files along with a publication, allowing independent replication even without relying
690 on *slendr* itself.

691 Moving forward, we expect that the *slendr* framework will become a useful tool [to produce](#)
692 [ground-truth data](#) for comparing and benchmarking inference methods for modeling spatial
693 genomic processes (Peter and Slatkin, 2013; Petkova, Novembre and Stephens, 2016; Marcus
694 *et al.*, 2021; Muktupavela *et al.*, 2021), [as well as for](#) the development of new approaches to
695 spatial problems in population genomics. There is great potential for deploying *slendr* in
696 simulation-based inference methods, [like Approximate Bayesian Computation \(ABC\)](#) (Beaumont,
697 [Zhang and Balding, 2002; Csilléry *et al.*, 2010](#)), thanks to its tight integration with the rest of the
698 R modeling landscape. [A major challenge in ABC is the significant amount of coding needed to](#)
699 [program simulations of demographic history and integrate them with software for computing](#)
700 [population genetic statistics. *slendr* can program complex models and compute relevant statistics](#)
701 [using its tree-sequence interface with a relatively small amount of code, all within a single R](#)
702 [workflow. Furthermore, although *slendr* does not currently include features for implicit, automated](#)
703 [parallelism \(an important aspect of computation-heavy modeling approaches such as ABC\), users](#)
704 [can rely on numerous R packages providing a wide range of parallelization techniques](#)
705 [\(Eddelbuettel, 2021\).](#)

706 Nonetheless, inference of spatial dynamics from genetic data remains an open research
707 problem with many potential pitfalls, and we strongly caution users to avoid overinterpretation.
708 For instance, *slendr* models retain a notion of discretely delineated populations, but even a
709 reasonable fit of such a model to real data does not erase the reality that such groupings are
710 rarely, if ever, as stable and cleanly distinguished as in idealized models. Indeed, confounding
711 the simple models used in population genetics with reality can be actively harmful (Coop, 2022;
712 [Khan *et al.*, 2022](#)). Furthermore, population genetic modeling in general is notoriously challenging
713 due to the many parameters involved ([Gravel *et al.*, 2011](#); [Pickrell and Pritchard, 2012](#); [Kamm *et*](#)
714 [al.](#), 2020), [In this respect, advanced, explicitly spatial models of the kind unlocked by *slendr*](#)
715 [present an even bigger challenge. For instance, how can we best do model comparison, and](#)
716 [among what set of models? What would constitute a good "null hypothesis" when modeling](#)
717 [potentially complex spatial population dynamics? Furthermore, even relatively simple models can](#)
718 [be ill-posed or even nonidentifiable: many combinations of spatial parameters \(such as individual](#)
719 [dispersal or mating distances\) may give rise to similar genetic patterns. Every demographic](#)
720 [inference study makes assumptions about the process which generated the data, sometimes](#)
721 [explicitly and sometimes implicitly, and awareness of these assumptions is vital for interpretation](#)
722 [of the results \(Loog, 2021\). We hope that the ease with which *slendr* allows one to explore the](#)
723 [impact of spatio-temporal parameters on population dynamics—and the fact that *slendr* forces](#)
724 [the researcher to state those parameters explicitly—will help guide researchers in establishing](#)
725 [guidelines for good practice, to delineate the limits of what can be learned and, consequently,](#)
726 [avoid overinterpretation \(or misinterpretation\) of such parameters.](#)

727 In its current version ([v0.5.0 as available on the CRAN repository](#)), *slendr*'s spatial
728 simulation maps are limited to landscapes that exhibit binary habitability—i.e., any given location
729 either is or is not habitable by individuals. [A more ecologically realistic simulation could allow for](#)
730 [varying degrees of habitability at different locations, which would affect the size of the simulated](#)

Deleted: .

Deleted: It will also enable

Deleted: especially

Deleted: , including methods such as Approximate Bayesian Computation (Beaumont, Zhang and Balding, 2002)

Deleted: However

Deleted: as

Deleted: in inference (population divergence times, times and magnitudes of gene-flow events and population size changes, etc.)

Deleted: , care will have to be taken to ameliorate the curse of dimensionality

Deleted: , without any intermediate degrees of habitability

Deleted: Future extensions of the *slendr* framework could include the incorporation of fine-scaled geographic raster maps which would affect local carrying capacities and influence population movements by introducing a cost of occupying a particular location.

751 [population. Future extensions of the *slendr* framework could include the incorporation of fine-](#)
752 [scaled geographic maps storing individual habitability values for each pixel of the raster, allowing](#)
753 [for dynamic changes of such maps over time. This would effectively make the size of the](#)
754 [population an emergent consequence of the habitability metric aggregated across the map. This](#)
755 extension would require significant changes to the *slendr* back-end code, moving to modeling
756 population densities per unit of landscape area using non-Wright–Fisher dynamics, but the
757 necessary software building blocks are already supported by SLiM [and examples of these types](#)
758 [of simulations are discussed in the SLiM manual](#) (Haller and Messer, 2022). A recently published
759 Python module *Geonomics* provides an interface for simulating genetic data on arbitrary
760 landscape rasters (Terasaki Hart, Bishop and Wang, 2021). Implementing such functionality in
761 *slendr* would have the advantage of using a much more efficient SLiM simulation engine and a
762 greater ease of use due to *slendr*'s emphasis on visually-focused interactive model design in R.
763 The main challenge would therefore lie in making sure that the additional complexity involved in
764 making the *slendr*'s SLiM back end more flexible does not compromise the current simplicity of
765 its declarative interface. The benefits of this extension would be numerous, including for genomic
766 forecasting and predicting species ranges in the face of climate change and ecological breakdown
767 (Fitzpatrick and Keller, 2015; Exposito-Alonso *et al.*, 2019; Theodoridis *et al.*, 2020), and for
768 constructing models of species distribution dynamics in the ancient past (Wang *et al.*, 2021).
769 [Implementation of this extension of *slendr* is still in the planning stages, in collaboration with the](#)
770 [community on the project's GitHub page.](#)

771 At the moment, *slendr* can only produce genome sequences from a single species
772 [\(although with an arbitrary number and spatial arrangement of population groups\)](#) due to the
773 restrictions imposed by its simulation back end. However, many types of genomic resources
774 distributed across space and time are represented by fragmentary mixtures of genomes from
775 multiple species. This includes ancient microbiomes from human remains (Rasmussen *et al.*,
776 2015), sedimentary DNA from permafrost, caves, or lake and marine cores (Willerslev *et al.*, 2003;
777 Parducci *et al.*, 2017; Armbrrecht *et al.*, 2019; Vernet *et al.*, 2021), and environmental DNA from
778 water, soil, or air samples (Taberlet *et al.*, 2012; Stat *et al.*, 2017; Lynggaard *et al.*, 2022). Recent
779 developments in SLiM would allow *slendr* to perform multi-species simulations, which would
780 facilitate ecological modeling of species distributions (Fordham *et al.*, 2021) or of past epidemics
781 (Duchene *et al.*, 2020) from a fully genomic perspective.

782 [Finally, at the time of writing, *slendr* models are limited to neutral simulations, and this](#)
783 [restriction applies even to simulations performed via its SLiM back end. In particular, *slendr* does](#)
784 [not currently provide built-in support for specifying mutation types, genomic element types,](#)
785 [recombination maps, or custom SLiM callbacks. Providing an R equivalent for SLiM's complete](#)
786 [functionality would be a daunting task of limited utility, and would substantially complicate *slendr*'s](#)
787 [intuitive R syntax for encoding demographic models \(Figure 1\). An attractive alternative for](#)
788 [supporting more advanced, customized models could be to retain the behavior of *slendr* described](#)
789 [in this manuscript as the default, but provide the possibility of overriding different aspects of this](#)
790 [behavior by injecting user-defined SLiM snippets at appropriate locations in *slendr*'s SLiM back-](#)
791 [end code. We are exploring this possibility for future versions of the software.](#)

792 Ultimately, we hope that our new simulation framework will help generate new ideas about
793 the insights that can be gleaned from the rich spatio-temporal information hidden within DNA
794 sequences. Furthermore, we aspire to help budding researchers in population genetics get started

Deleted: populations of

796 with simulations and build their intuition about population genetic concepts by developing models
797 using more traditional non-spatial methods and statistics, and we believe that *slendr* could be a
798 useful tool for teaching population genetics to students. We hope that by easily generating and
799 visualizing genomic models on real landscapes, we can spark new ways of thinking about how
800 organisms evolve (Bradburd and Ralph, 2019) and enable clearer discussions about the
801 fundamental interconnectedness of genomes across space and time (Mathieson and Scally,
802 2020).

803 Acknowledgements

804 We thank Moisés Coll Macià, Mariadaria Kathrine Ianni-Ravn, and Emma Prantoni for testing and
805 feedback on early versions of *slendr*, and the members of the Racimo group for valuable
806 comments on the design of the software and this manuscript. FR was supported by a Villum
807 Young Investigator Grant (project no. 00025300), COREX ERC Synergy grant (ID 951385),
808 a Novo Nordisk Fonden Data Science Ascending Investigator Award (NNF22OC0076816)
809 and a Sapere Aude grant (2064-00026B) from Danmarks Frie Forskningsfond. MP was
810 supported by a Lundbeck Foundation grant (R302-2018-2155) and a Novo Nordisk
811 Foundation grant (NNF18SA0035006) given to the GeoGenetics Centre. PR was supported
812 by NIH award R01HG010774.

813 References

- 814 1000 Genomes Project (2010) 'A map of human genome variation from population-scale
815 sequencing', *Nature*, 467(7319), pp. 1061–1073.
- 816 Al-Asadi, H. *et al.* (2019) 'Estimating recent migration and population-size surfaces', *PLoS*
817 *genetics*, 15(1), p. e1007908.
- 818 Allentoft, M.E. *et al.* (2015) 'Population genomics of Bronze Age Eurasia', *Nature*, 522(7555),
819 pp. 167–172.
- 820 Armbrrecht, L.H. *et al.* (2019) 'Ancient DNA from marine sediments: Precautions and
821 considerations for seafloor coring, sample handling and data generation', *Earth-Science*
822 *Reviews*, 196, p. 102887.
- 823 Barton, N. (1979) 'Gene flow past a cline', *Heredity*, 43(3), pp. 333–339.
- 824 Barton, N., Depaulis, F. and Etheridge, A.M. (2002) 'Neutral evolution in spatially continuous
825 populations', *Theoretical population biology*, 61(1), pp. 31–48.
- 826 Barton, N., Etheridge, A.M. and Véber, A. (2013) 'Modelling evolution in a spatial continuum',
827 *Journal of Statistical Mechanics: Theory and Experiment*, 2013(01), p. P01002.
- 828 Barton, N., Etheridge, A. and Véber, A. (2010) 'A New Model for Evolution in a Spatial
829 Continuum', *Electronic Journal of Probability*, 15.
- 830 Baumdicker, F. *et al.* (2022) 'Efficient ancestry and mutation simulation with msprime 1.0',
831 *Genetics*, 220(3). Available at: <https://doi.org/10.1093/genetics/iyab229>.

- 832 Beaumont, M.A., Zhang, W. and Balding, D.J. (2002) 'Approximate Bayesian computation in
833 population genetics', *Genetics*, 162(4), pp. 2025–2035.
- 834 Beerli, P. and Felsenstein, J. (2001) 'Maximum likelihood estimation of a migration matrix and
835 effective population sizes in n subpopulations by using a coalescent approach', *Proceedings of
836 the National Academy of Sciences of the United States of America*, 98(8), pp. 4563–4568.
- 837 Bergström, A. *et al.* (2020) 'Origins and genetic legacy of prehistoric dogs', *Science*, 370(6516),
838 pp. 557–564.
- 839 Bradburd, G.S., Coop, G.M. and Ralph, P.L. (2018) 'Inferring Continuous and Discrete
840 Population Genetic Structure Across Space', *Genetics*, 210(1), pp. 33–52.
- 841 Bradburd, G.S. and Ralph, P.L. (2019) 'Spatial population genetics: It's about time', *Annual
842 review of ecology, evolution, and systematics*, 50(1), pp. 427–449.
- 843 Chambers, J.M. (2020) 'S, R, and data science', *Proceedings of the ACM on Programming
844 Languages*, pp. 1–17. Available at: <https://doi.org/10.1145/3386334>.
- 845 Chang, W. *et al.* (2021) *shiny: Web Application Framework for R*. Available at: [https://CRAN.R-
846 project.org/package=shiny](https://CRAN.R-project.org/package=shiny).
- 847 Coop, G. (2022) 'Genetic similarity versus genetic ancestry groups as sample descriptors in
848 human genetics', *arXiv [q-bio.PE]*. Available at: <http://arxiv.org/abs/2207.11595>.
- 849 Crabtree, S.A. *et al.* (2021) 'Landscape rules predict optimal superhighways for the first
850 peopling of Sahul', *Nature human behaviour*, 5(10), pp. 1303–1313.
- 851 Csilléry, K. *et al.* (2010) 'Approximate Bayesian Computation (ABC) in practice', *Trends in
852 ecology & evolution*, 25(7), pp. 410–418.
- 853 Currat, M. *et al.* (2019) 'SPLATCHE3: simulation of serial genetic data under spatially explicit
854 evolutionary scenarios including long-distance dispersal', *Bioinformatics*, 35(21), pp. 4480–
855 4483.
- 856 Currat, M. and Excoffier, L. (2011) 'Strong reproductive isolation between humans and
857 Neanderthals inferred from observed patterns of introgression', *Proceedings of the National
858 Academy of Sciences*, 108(37), pp. 15129–15134.
- 859 Currat, M., Ray, N. and Excoffier, L. (2004) 'splatche: a program to simulate genetic diversity
860 taking into account environmental heterogeneity', *Molecular ecology notes*, 4(1), pp. 139–142.
- 861 Danecek, P. *et al.* (2011) 'The variant call format and VCFtools', *Bioinformatics*, 27(15), pp.
862 2156–2158.
- 863 Delsler, P.M. *et al.* (2021) 'Climate and mountains shaped human ancestral genetic lineages',
864 *bioRxiv*. Available at: <https://doi.org/10.1101/2021.07.13.452067>.
- 865 Duchene, S. *et al.* (2020) 'Temporal signal and the phylodynamic threshold of SARS-CoV-2',
866 *Virus evolution*, 6(2), p. veaa061.
- 867 Duforet-Frebourg, N. and Blum, M.G.B. (2014) 'Nonstationary patterns of isolation-by-distance:
868 inferring measures of local genetic differentiation with Bayesian kriging', *Evolution; international*

Commented [kk1]: Citation added during revisions.

- 869 *journal of organic evolution*, 68(4), pp. 1110–1123.
- 870 Eddelbuettel, D. (2021) 'Parallel computing with R: A brief review', *Wiley interdisciplinary*
871 *reviews. Computational statistics*, 13(2). Available at: <https://doi.org/10.1002/wics.1515>.
- 872 Ewing, G. and Hermisson, J. (2010) 'MSMS: a coalescent simulation program including
873 recombination, demographic structure and selection at a single locus', *Bioinformatics*, 26(16),
874 pp. 2064–2065.
- 875 Exposito-Alonso, M. *et al.* (2019) 'Natural selection on the *Arabidopsis thaliana* genome in
876 present and future climates', *Nature*, 573(7772), pp. 126–129.
- 877 Felsenstein, J. (1975) 'A Pain in the Torus: Some Difficulties with Models of Isolation by
878 Distance', *The American naturalist*, 109(967), pp. 359–368.
- 879 Feuerborn, T.R. *et al.* (2021) 'Modern Siberian dog ancestry was shaped by several thousand
880 years of Eurasian-wide trade and human dispersal', *Proceedings of the National Academy of*
881 *Sciences of the United States of America*, 118(39). Available at:
882 <https://doi.org/10.1073/pnas.2100338118>.
- 883 Fisher, R.A. (1937) 'The wave of advance of advantageous genes', *Annals of eugenics*, 7(4),
884 pp. 355–369.
- 885 Fitzpatrick, M.C. and Keller, S.R. (2015) 'Ecological genomics meets community-level modelling
886 of biodiversity: mapping the genomic landscape of current and future environmental adaptation',
887 *Ecology letters*, 18(1), pp. 1–16.
- 888 Fordham, D.A. *et al.* (2021) 'poems: R package for simulating species' range dynamics using
889 pattern-oriented validation', *Methods in ecology and evolution / British Ecological Society*,
890 12(12), pp. 2364–2371.
- 891 Frachetti, M.D. *et al.* (2017) 'Nomadic ecology shaped the highland geography of Asia's Silk
892 Roads', *Nature*, 543(7644), pp. 193–198.
- 893 Fu, Q. *et al.* (2016) 'The genetic history of Ice Age Europe', *Nature*, 534(7606), pp. 200–205.
- 894 Gravel, S. *et al.* (2011) 'Demographic history and rare allele sharing among human populations',
895 *Proceedings of the National Academy of Sciences*, 108(29), pp. 11983–11988.
- 896 Green, R.E. *et al.* (2010) 'A draft sequence of the Neandertal genome', *Science*, 328(5979), pp.
897 710–722.
- 898 Guillot, G. *et al.* (2009) 'Statistical methods in spatial genetics', *Molecular ecology*, 18(23), pp.
899 4734–4756.
- 900 Haak, W. *et al.* (2015) 'Massive migration from the steppe was a source for Indo-European
901 languages in Europe', *Nature*, 522(7555), pp. 207–211.
- 902 Haller, B.C. *et al.* (2019) 'Tree-sequence recording in SLiM opens new horizons for forward-time
903 simulation of whole genomes', *Molecular ecology resources*, 19(2), pp. 552–566.
- 904 Haller, B.C. and Messer, P.W. (2017) 'SLiM 2: Flexible, Interactive Forward Genetic
905 Simulations', *Molecular biology and evolution*, 34(1), pp. 230–240.

- 906 Haller, B.C. and Messer, P.W. (2019) 'SLiM 3: Forward Genetic Simulations Beyond the Wright-
907 Fisher Model', *Molecular biology and evolution*, 36(3), pp. 632–637.
- 908 Haller, B.C. and Messer, P.W. (2022) 'SLiM 4: Multispecies eco-evolutionary modeling', *The
909 American naturalist* [Preprint]. Available at: <https://doi.org/10.1086/723601>.
- 910 Hanks, E.M. and Hooten, M.B. (2013) 'Circuit theory and model-based inference for landscape
911 connectivity', *Journal of the American Statistical Association*, 108(501), pp. 22–33.
- 912 Hudson, R.R. (2002) 'Generating samples under a Wright-Fisher neutral model of genetic
913 variation', *Bioinformatics*, 18(2), pp. 337–338.
- 914 Kamm, J. *et al.* (2020) 'Efficiently inferring the demographic history of many populations with
915 allele count data', *Journal of the American Statistical Association*, 115(531), pp. 1472–1487.
- 916 Kelleher, J. *et al.* (2018) 'Efficient pedigree recording for fast population genetics simulation',
917 *PLoS computational biology*, 14(11), p. e1006581.
- 918 Kelleher, J. *et al.* (2019) 'Inferring whole-genome histories in large population datasets', *Nature
919 genetics*, 51(9), pp. 1330–1338.
- 920 Kelleher, J., Etheridge, A.M. and Barton, N.H. (2014) 'Coalescent simulation in continuous
921 space: algorithms for large neighbourhood size', *Theoretical population biology*, 95, pp. 13–23.
- 922 Kelleher, J., Etheridge, A.M. and McVean, G. (2016) 'Efficient Coalescent Simulation and
923 Genealogical Analysis for Large Sample Sizes', *PLoS computational biology*, 12(5), p.
924 e1004842.
- 925 Khan, A.T. *et al.* (2022) 'Recommendations on the use and reporting of race, ethnicity, and
926 ancestry in genetic research: Experiences from the NHLBI TOPMed program', *Cell genomics*,
927 2(8). Available at: <https://doi.org/10.1016/j.xgen.2022.100155>.
- 928 Kimura, M. (1953) "'Stepping Stone" model of population', *Annual report of the National Institute
929 of Genetics*, 3, pp. 62–63.
- 930 Kimura, M. and Weiss, G.H. (1964) 'The Stepping Stone Model of Population Structure and the
931 Decrease of Genetic Correlation with Distance', *Genetics*, 49(4), pp. 561–576.
- 932 Lazaridis, I. *et al.* (2014) 'Ancient human genomes suggest three ancestral populations for
933 present-day Europeans', *Nature*, 513(7518), pp. 409–413.
- 934 Levene, H. (1953) 'Genetic Equilibrium When More Than One Ecological Niche is Available',
935 *The American naturalist*, 87(836), pp. 331–333.
- 936 Librado, P. *et al.* (2021) 'The origins and spread of domestic horses from the Western Eurasian
937 steppes', *Nature*, pp. 1–7.
- 938 Liu, H. *et al.* (2006) 'A geographically explicit genetic model of worldwide human-settlement
939 history', *American journal of human genetics*, 79(2), pp. 230–237.
- 940 Loog, L. *et al.* (2017) 'Estimating mobility using sparse data: Application to human genetic
941 variation', *Proceedings of the National Academy of Sciences*, 114(46), pp. 12213–12218.

Commented [kk2]: Citation added during revisions.

- 942 Loog, L. (2021) 'Sometimes hidden but always there: the assumptions underlying genetic
943 inference of demographic histories', *Philosophical transactions of the Royal Society of London.*
944 *Series B, Biological sciences*, 376(1816), p. 20190719.
- 945 Lovelace, R., Nowosad, J. and Muenchow, J. (2019) *Geocomputation with R (Chapman &*
946 *Hall/CRC The R Series)*. 1st edn. Routledge.
- 947 Lynggaard, C. *et al.* (2022) 'Airborne environmental DNA for terrestrial vertebrate community
948 monitoring', *Current biology: CB*, 32(3), pp. 701–707.e5.
- 949 Malécot, G. (1951) 'Un traitement stochastique des problèmes linéaires en génétique de
950 population', *Ann. Univ. Lyon. Sci. Sec.*, 14, pp. 79–117.
- 951 Mallick, S. *et al.* (2016) 'The Simons Genome Diversity Project: 300 genomes from 142 diverse
952 populations', *Nature*, 538(7624), pp. 201–206.
- 953 Marcus, J. *et al.* (2021) 'Fast and flexible estimation of effective migration surfaces', *eLife*, 10.
954 Available at: <https://doi.org/10.7554/eLife.61927>.
- 955 Mathieson, I. and Scally, A. (2020) 'What is ancestry', *PLoS genetics*, 16(3), p. e1008624.
- 956 McRae, B.H. (2006) 'Isolation by resistance', *Evolution; international journal of organic*
957 *evolution*, 60(8), pp. 1551–1561.
- 958 Meyer, M. *et al.* (2017) 'Palaeogenomes of Eurasian straight-tusked elephants challenge the
959 current view of elephant evolution', *eLife*, 6, p. e25413.
- 960 Muktupavela, R.A. *et al.* (2022) 'Modeling the spatiotemporal spread of beneficial alleles using
961 ancient genomes', *eLife*, 11, p. e73767.
- 962 Nagylaki, T. (1976) 'The relation between distant individuals in geographically structured
963 populations', *Mathematical biosciences*, 28(1), pp. 73–80.
- 964 Osmond, M.M. and Coop, G. (2021) 'Estimating dispersal rates and locating genetic ancestors
965 with genome-wide genealogies', *bioRxiv*. Available at:
966 <https://doi.org/10.1101/2021.07.13.452277>.
- 967 Palkopoulou, E. *et al.* (2018) 'A comprehensive genomic history of extinct and living elephants',
968 *Proceedings of the National Academy of Sciences of the United States of America*, 115(11), pp.
969 E2566–E2574.
- 970 Paradis, E. (2011) *Analysis of Phylogenetics and Evolution with R (Use R!)*. 2nd edn. Springer.
- 971 Paradis, E. and Schliep, K. (2019) 'ape 5.0: an environment for modern phylogenetics and
972 evolutionary analyses in R', *Bioinformatics*. Edited by R. Schwartz, 35(3), pp. 526–528.
- 973 Parducci, L. *et al.* (2017) 'Ancient plant DNA in lake sediments', *The New phytologist*, 214(3),
974 pp. 924–942.
- 975 Patterson, N. *et al.* (2012) 'Ancient admixture in human history', *Genetics*, 192(3), pp. 1065–
976 1093.
- 977 Pebesma, E. (2018) 'Simple features for R: Standardized support for spatial vector data', *The R*

Commented [kk3]: Citation added during revisions.

- 978 *journal*, 10(1), p. 439.
- 979 Peter, B.M. and Slatkin, M. (2013) 'Detecting range expansions from genetic data', *Evolution; international journal of organic evolution*, 67(11), pp. 3274–3289.
- 980
- 981 Petkova, D., Novembre, J. and Stephens, M. (2016) 'Visualizing spatial population structure with
982 estimated effective migration surfaces', *Nature genetics*, 48(1), pp. 94–100.
- 983 Pickrell, J.K. and Pritchard, J.K. (2012) 'Inference of Population Splits and Mixtures from
984 Genome-Wide Allele Frequency Data', *PLoS genetics*. Edited by H. Tang, 8(11), p. e1002967.
- 985 Pickrell, J.K. and Reich, D. (2014) 'Toward a new history and geography of human genes
986 informed by ancient DNA', *Trends in genetics: TIG*, 30(9), pp. 377–389.
- 987 Purcell, S. *et al.* (2007) 'PLINK: a tool set for whole-genome association and population-based
988 linkage analyses', *American journal of human genetics*, 81(3), pp. 559–575.
- 989 Racimo, F. *et al.* (2020) 'The spatiotemporal spread of human migrations during the European
990 Holocene', *Proceedings of the National Academy of Sciences of the United States of America*,
991 117(16), pp. 8989–9000.
- 992 Ralph, P. and Coop, G. (2013) 'The Geography of Recent Genetic Ancestry across Europe',
993 *PLoS biology*. Edited by C. Tyler-Smith, 11(5), p. e1001555.
- 994 Ralph, P., Thornton, K. and Kelleher, J. (2020) 'Efficiently Summarizing Relationships in Large
995 Samples: A General Duality Between Statistics of Genealogies and Genomes', *Genetics*,
996 215(3), pp. 779–797.
- 997 Rasmussen, M. *et al.* (2010) 'Ancient human genome sequence of an extinct Palaeo-Eskimo',
998 *Nature*, 463(7282), pp. 757–762.
- 999 Rasmussen, S. *et al.* (2015) 'Early divergent strains of *Yersinia pestis* in Eurasia 5,000 years
1000 ago', *Cell*, 163(3), pp. 571–582.
- 1001 R Core Team (2021) 'R: A Language and Environment for Statistical Computing'. Vienna,
1002 Austria: R Foundation for Statistical Computing. Available at: <https://www.R-project.org/>.
- 1003 Ringbauer, H. *et al.* (2018) 'Estimating Barriers to Gene Flow from Distorted Isolation-by-
1004 Distance Patterns', *Genetics*, 208(3), pp. 1231–1245.
- 1005 Rousset, F. (1997) 'Genetic differentiation and estimation of gene flow from F-statistics under
1006 isolation by distance', *Genetics*, 145(4), pp. 1219–1228.
- 1007 Safner, T. *et al.* (2011) 'Comparison of Bayesian clustering and edge detection methods for
1008 inferring boundaries in landscape genetics', *International journal of molecular sciences*, 12(2),
1009 pp. 865–889.
- 1010 Sandve, G.K. *et al.* (2013) 'Ten simple rules for reproducible computational research', *PLoS
1011 computational biology*, 9(10), p. e1003285.
- 1012 Slatkin, M. (1973) 'GENE FLOW AND SELECTION IN A CLINE', *Genetics*, 75(4), pp. 733–756.
- 1013 Slatkin, M. and Excoffier, L. (2012) 'Serial founder effects during range expansion: a spatial

- 1014 analog of genetic drift', *Genetics*, 191(1), pp. 171–181.
- 1015 Slatkin, M. and Racimo, F. (2016) 'Ancient DNA and human history', *Proceedings of the*
1016 *National Academy of Sciences*, 113(23), pp. 6380–6387.
- 1017 Speidel, L. *et al.* (2019) 'A method for genome-wide genealogy estimation for thousands of
1018 samples', *Nature genetics*, 51(9), pp. 1321–1329.
- 1019 Staab, P.R. *et al.* (2015) 'scrm: efficiently simulating long sequences using the approximated
1020 coalescent with recombination', *Bioinformatics*, 31(10), pp. 1680–1682.
- 1021 Stat, M. *et al.* (2017) 'Ecosystem biomonitoring with eDNA: metabarcoding across the tree of life
1022 in a tropical marine environment', *Scientific reports*, 7(1), p. 12240.
- 1023 Taberlet, P. *et al.* (2012) 'Towards next-generation biodiversity assessment using DNA
1024 metabarcoding', *Molecular ecology*, 21(8), pp. 2045–2050.
- 1025 Terasaki Hart, D.E., Bishop, A.P. and Wang, I.J. (2021) 'Geonomics: Forward-Time, Spatially
1026 Explicit, and Arbitrarily Complex Landscape Genomic Simulations', *Molecular biology and*
1027 *evolution*, 38(10), pp. 4634–4646.
- 1028 Theodoridis, S. *et al.* (2020) 'Evolutionary history and past climate change shape the distribution
1029 of genetic diversity in terrestrial mammals', *Nature communications*, 11(1), p. 2557.
- 1030 Vernot, B. *et al.* (2021) 'Unearthing Neanderthal population history using nuclear and
1031 mitochondrial DNA from cave sediments', *Science*, 372(6542). Available at:
1032 <https://doi.org/10.1126/science.abf1667>.
- 1033 Wang, Y. *et al.* (2021) 'Late Quaternary dynamics of Arctic biota from ancient environmental
1034 genomics', *Nature*, 600(7887), pp. 86–92.
- 1035 Wickham, H. *et al.* (2019) 'Welcome to the tidyverse', *Journal of Open Source Software*, p.
1036 1686. Available at: <https://doi.org/10.21105/joss.01686>.
- 1037 Wickham, H. and Grolemund, G. (2016) *R for Data Science*. "O'Reilly Media, Inc.", p. 520.
- 1038 Willerslev, E. *et al.* (2003) 'Diverse plant and animal genetic records from Holocene and
1039 Pleistocene sediments', *Science*, 300(5620), pp. 791–795.
- 1040 Wohns, A.W. *et al.* (2022) 'A unified genealogy of modern and ancient genomes', *Science*,
1041 375(6583), p. eabi8264.
- 1042 Wright, S. (1943) 'Isolation by Distance', *Genetics*, 28(2), pp. 114–138.
- 1043 Yu, G. *et al.* (2017) 'ggtree : An R package for visualization and annotation of phylogenetic trees
1044 with their covariates and other associated data', *Methods in ecology and evolution / British*
1045 *Ecological Society*, 8(1), pp. 28–36.