

¹ Estimating dispersal rates and locating genetic
² ancestors with genome-wide genealogies

³ Matthew M Osmond¹ and Graham Coop²

⁴ ¹Department of Ecology & Evolutionary Biology, University of
⁵ Toronto

⁶ ²Department of Evolution & Ecology and Center for
⁷ Population Biology, University of California - Davis

⁸ **Abstract**

⁹ Spatial patterns in genetic diversity are shaped by individuals dispersing
¹⁰ from their parents and larger-scale population movements. It has long been
¹¹ appreciated that these patterns of movement shape the underlying genealo-
¹² gies along the genome leading to geographic patterns of isolation by distance
¹³ in contemporary population genetic data. However, extracting the enormous
¹⁴ amount of information contained in genealogies along recombining sequences
¹⁵ has, up till recently, not been computational feasible. Here we capitalize
¹⁶ on important recent advances in gene-genealogy reconstruction and develop
¹⁷ methods to use thousands of trees to estimate time-varying per-generation
¹⁸ dispersal rates and to locate the genetic ancestors of a sample back through
¹⁹ time. We take a likelihood approach using a simple approximate model
²⁰ (branching Brownian motion) as our prior distribution of spatial genealogies.
²¹ After testing our method with simulations we apply it to the 1001 Genomes
²² dataset of over one thousand *Arabidopsis thaliana* genomes sampled across
²³ a wide geographic extent. We detect a very high dispersal rate in the recent
²⁴ past, especially longitudinally, and use inferred ancestor locations to visualize
²⁵ many examples of recent long-distance dispersal and admixture. We also use
²⁶ inferred ancestor locations to identify the origin and ancestry of the North

27 American expansion and to depict alternative geographic ancestries stemming from multiple glacial refugia. Our method highlights the huge amount
28 of information about past dispersal events and population movements contained in genome-wide genealogies.
29

31 Introduction

32 Patterns of genetic diversity are shaped by the movements of individuals, as
33 individuals move their alleles around the landscape as they disperse. Patterns
34 of individual movement reflect individual-level dispersal; children move away
35 from their parents village and dandelion seeds blow in the wind. These
36 patterns also reflect large-scale movements of populations. For example, in
37 the past decade we have learnt about the large-scale movement of different
38 peoples across the world from ancient DNA ([Slatkin and Racimo, 2016](#); [Reich, 2018](#)). Such large scale movements of individuals also occur in other species
40 during biological invasions or with the retreat and expansion of populations
41 in and out of glacial refugia, tracking the waxing and waning of the ice ages
42 ([Hewitt, 2000](#)).

43 An individual's set of genealogical ancestors expands rapidly geographically
44 back through time in sexually reproducing organisms ([Kelleher et al., 2016a](#); [Coop, 2017](#)). Due to limited recombination each generation, more
46 than a few tens of generations back an individual's genetic ancestors represent
47 only a tiny sample of their vast number of genealogical ancestors ([Donnelly, 1983](#); [Coop, 2013](#)). Yet the geographic locations of genetic ancestors
49 still represent an incredibly rich source of information on population history
50 ([Bradburd and Ralph, 2019](#)). We can hope to learn about the geography of
51 genetic ancestors because individuals who are geographically close are often
52 genetically more similar across their genomes; their ancestral lineages have
53 only dispersed for a relatively short time and distance since they last shared
54 a geographically-close common ancestor. This pattern is termed isolation
55 by distance. These ideas about the effects of geography and genealogy have
56 underlain our understanding of spatial population genetics since its inception
57 ([Wright, 1943](#); [Malécot, 1948](#)). Under coalescent models, lineages move
58 spatially, as a Brownian motion if dispersal is random and local, splitting to
59 give rise to descendent lineages till we reach the present day. Such models
60 underlie inferences based on increasing allele frequency differentiation (such
61 as F_{ST}) with geographic distance ([Rousset, 1997](#)) and the drop-off in the

sharing of long blocks of genome shared identical by descent among pairs of individuals (Ralph and Coop, 2013; Ringbauer et al., 2017). These models also are the basis of methods that seek violations of isolation-by-distance (Wang and Bradburd, 2014).

While spatial genealogies have proven incredibly useful for theoretical tools and intuition, with few exceptions they have not proven useful for inferences because we have not been able to construct these genealogies along recombining sequences. In non-recombining chromosomes (e.g., mtDNA and Y), constructed genealogies have successfully been used to understand patterns of dispersal and spatial spread (Avise, 2009). However, these spatial genealogy inferences are necessarily limited as a single genealogy holds only limited information about the history of populations in a recombining species (Barton and Wilson, 1995). Phylogenetic approaches to geography ('phylogeography'; Knowles, 2009) have been more widely and successfully applied to pathogens to track the spatial spread of epidemics, such as SARS-CoV-2 (Martin et al., 2021), but such approaches have yet to be applied to the thousands of genealogical trees that exist in sexual populations.

Here we capitalize on the recent ability to infer a sequence of genealogies, with branch lengths, across recombining genomes (Rasmussen et al., 2014; Speidel et al., 2019; Wohns et al., 2021). Equipped with this information, we develop a method that uses a sequence of trees to estimate dispersal rates and locate genetic ancestors in continuous, 2-dimensional space, under the assumption of Brownian motion. Using thousands of approximately unlinked trees, we multiply likelihoods of dispersal rates across trees to get genome-wide estimates and use the sequence of trees to predict a cloud of ancestral locations as a way to visualize geographic ancestries. We first test our approach with simulations and then apply it to *Arabidopsis thaliana*, using over 1000 genomes with a wide geographic distribution (Alonso-Blanco et al., 2016) and a complex history (Fulgione and Hancock, 2018; Hsu et al., 2019).

Results

Overview of approach

We first give an overview of our approach, the major components of which are illustrated in Figure 1. See Materials and Methods for more details.

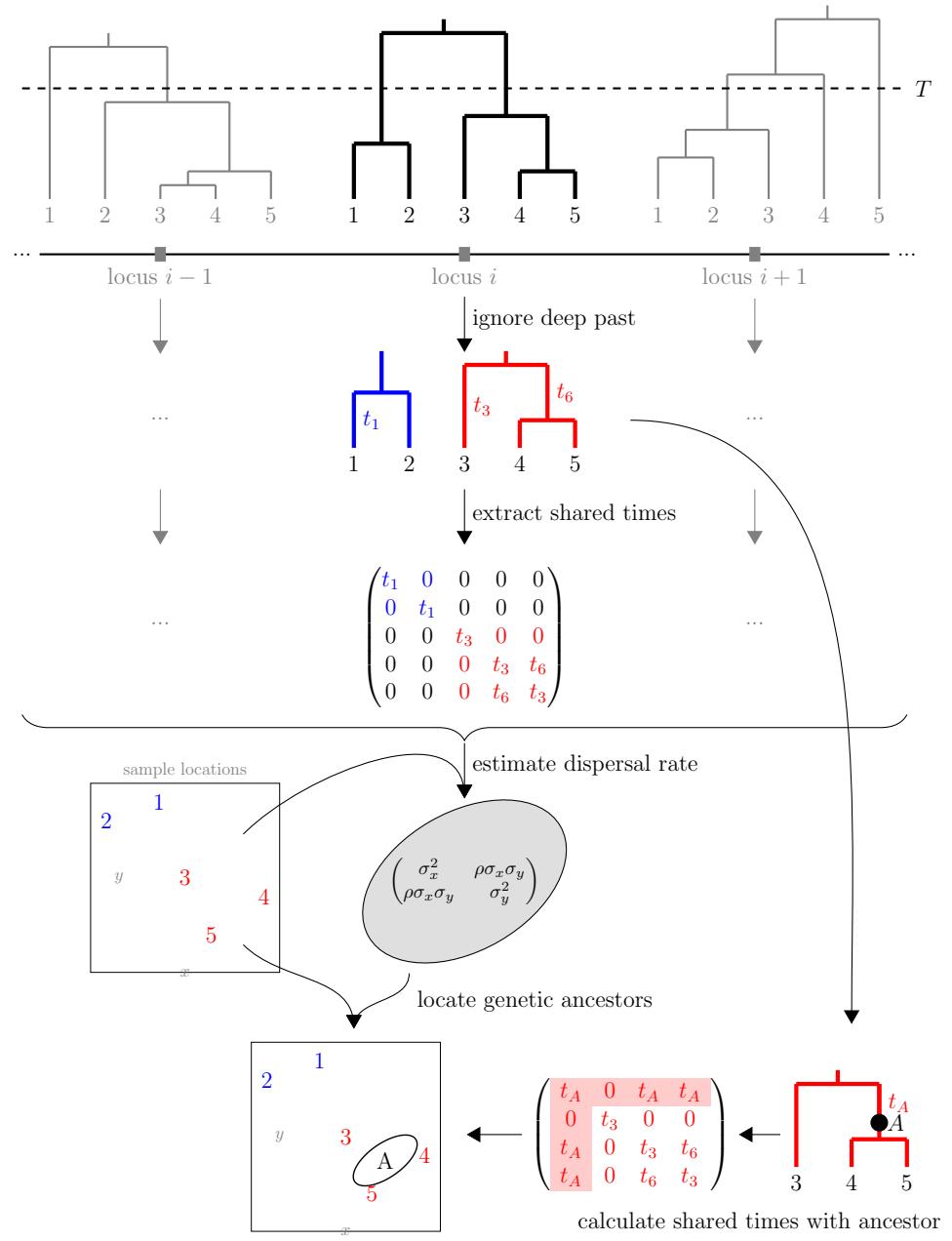


Figure 1: (Continued on the following page.)

Figure 1: **Conceptual overview of the approach.** From a sequence of trees covering the full genome, we downsample to trees at approximately unlinked loci. To avoid the influence of strongly non-Brownian dynamics at deeper times (e.g., glacial refugia, boundaries), we ignore times deeper than T , which divides each tree into multiple subtrees (here, blue and red subtrees at locus i). From these subtrees we extract the shared evolutionary times of each lineage with all others. In practice (but not shown here), we use multiple samples of the tree at a given locus, for importance sampling, and also extract the coalescence times for importance sample weights. Under Brownian motion, the shared times describe the covariance we expect to see in the locations of our samples, and so using the times and locations we can find maximum likelihood dispersal rate (a 2×2 covariance matrix). While we can estimate a dispersal rate at each locus, a strength of our approach is that we combine information across many loci, by multiplying likelihoods, to estimate genome-wide (or per-chromosome) dispersal rates. In practice (but not shown here), we estimate dispersal in multiple epochs, which allows for variation in dispersal rate over time and helps absorb non-Brownian movements further in the past. Finally, we locate a genetic ancestor at a particular locus (a point on a tree, here A) by first calculating the time this ancestor shares with each of the samples in its subtree, and then using this matrix and the dispersal rate to calculate the probability distribution of the ancestor's location conditioned on the sample locations. In practice (but not shown here), we calculate the locations of ancestors of a given sample at a given time at many loci, combining information across loci into a distribution of genome-wide ancestry across space.

96 The rate of dispersal, which determines the average distance between par-
97 ents and offspring, is a key parameter in ecology and evolution. To estimate
98 this parameter we assume that in each generation the displacement of an
99 offspring from its mother is normally distributed with a mean of 0 and co-
100 variance matrix Σ . The covariance matrix is determined by the standard
101 deviations σ_x and σ_y along the x (longitudinal) and y (latitudinal) axes, re-
102 spectively, as well as the correlation, ρ , in displacements between these two
103 axes. The average distance between mothers and offspring is then $\sqrt{2/\pi}\sigma_i$ in
104 each dimension. We refer to the covariance matrix Σ as the (per-generation)
105 dispersal rate.

106 Given this model, the path of a lineage from its ancestral location to the

107 present day location is described by a Brownian motion. Lineages covary
108 in their locations because of shared evolutionary histories – lineages with a
109 more recent common ancestor covary more. Given a tree at a locus we can
110 calculate the covariance matrix of shared evolutionary times and compute
111 the likelihood of the dispersal rate, Σ , which is normally distributed given
112 this covariance matrix (Equation (1)). At each locus we can estimate the
113 likelihood of the dispersal rate given the tree at that locus and, given we
114 sample loci far enough apart that they are essentially independent, multiply
115 likelihoods across loci to derive a genome-wide likelihood, and thus a genome-
116 wide maximum likelihood dispersal rate.

117 Under this same model we can also estimate the locations of genetic
118 ancestors at a locus. Any point along a tree at any locus is a genetic ancestor
119 of one or more current day samples. This ancestor's lineage has dispersed
120 away from the location of the most recent common ancestor of the sample,
121 and covaries with current day samples in their geographic location to the
122 extent that it shares times in the tree with them. Under this model, the
123 location of an ancestor is influenced by the locations of all samples in the
124 same tree, including those that are not direct descendants (cf. Wohns et al.,
125 2021). For example, in Figure 1 the ancestor's location is not be the midpoint
126 of its two descendants (samples 4 and 5); the ancestor's location is also pulled
127 towards sample 3 since the ancestor and sample 3 both arose from a common
128 ancestor. Conditioning on the sample locations, and given the shared times
129 and previously inferred dispersal rate, we can compute the probability the
130 ancestor was at any location, which again is a normal distribution (Equation
131 (23)). In contrast to dispersal, for ancestral locations we do not want to
132 multiply likelihoods across loci since the ancestors at distant loci are likely
133 distinct. Instead we calculate the maximum likelihood location of genetic
134 ancestors at each locus to get a cloud of likely ancestral locations, and use
135 these clouds to visualize the spatial spread of ancestry backwards through
136 time.

137 We estimate marginal trees along the genome using **Relate** (Speidel et al.,
138 2019). Relate infers a sequence of tree topologies and associated branch
139 lengths, and can return a set of posterior draws of the branch lengths on a
140 given tree. This posterior distribution of branch lengths is useful to us as the
141 shared times in the tree are key to the amount of time that individual lineages
142 have had to disperse away from one another and we wish include uncertainty
143 in the times into our method. Relate gives us a posterior distribution of
144 branch lengths that is estimated using a coalescent prior, which assumes a

145 panmictic population of varying population size (the size changes are esti-
146 mated as part of the method), where any two lineages are equally likely to
147 coalesce. This panmictic prior results in a bias in the coalescent times un-
148 der a spatial model, where geographically proximate samples are more likely
149 to coalesce. To correct for this bias we make use of importance sampling
150 to weight the samples of branch lengths at each locus. We then calculate
151 the weighted average likelihood over our draws of our sample of trees at a
152 locus (or loci), so that it is as if they were drawn from a prior of branching
153 Brownian motion ([Meligkotsidou and Fearnhead, 2007](#)). Branching Brow-
154 nian motion, also known as the Brownian-Yule process, is a simple model
155 of spatial genealogies in a continuous population; it does not describe the
156 full complexities of spatial models such as the spatial coalescent, but it pro-
157 vides an analytically tractable model and a reasonable approximation over
158 short-time scales ([Edwards, 1970](#); [Rannala and Yang, 1996](#); [Meligkotsidou](#)
159 and [Fearnhead, 2007](#); [Novembre and Slatkin, 2009](#)).

160 In practice we concentrate on the recent past history of our sample. For
161 our estimates of dispersal rates in particular we do not want to assume that
162 our model of Gaussian dispersal (and branching Brownian motion) holds deep
163 into the past history of the sample. This is because the long-term movement
164 of lineages is constrained by geographic barriers (e.g., oceans) and larger scale
165 population movements may erase geographic signals over deep time scales.
166 On theoretical level ignoring the deep past may also be justified because
167 in a finite habitat the locations of coalescence events further back in time
168 become independent of sampling locations as lineages have moved around
169 sufficiently ([Wilkins and Wakeley, 2002](#)). Thus we only use this geographic
170 model to some time point in the past (T), and at each locus we use the
171 covariance of shared branch lengths based on the set of subtrees formed by
172 cutting off the full tree T generations back.

173 To relax the assumption of a constant dispersal rate we extend our method
174 to estimate dispersal rates in multiple epochs. Under Brownian motion this
175 extension is fairly straight-forward as the covariance in sample locations is
176 simply the sum, across epochs, of epoch-specific covariances (the Kronecker
177 product of the shared times in a epoch and the dispersal rate in that epoch).
178 An added benefit of estimating dispersal in epochs is that the estimates in
179 more distant epochs can absorb some of the non-Brownian dynamics further
180 back in time, increasing the accuracy of estimates in more recent epochs.

181 **Simulations**

182 We first wanted to test the performance of our method in a situation where
183 the true answers were known. To do this we used a combination of spatially-
184 explicit forward-time simulations (Haller and Messer, 2019), coalescent sim-
185 ulations (Kelleher et al., 2016b), and tree-sequence tools (Haller et al., 2019;
186 Kelleher et al., 2019; Speidel et al., 2019) to compare our estimates of disper-
187 sal rates and ancestor locations with the truth (see Materials and Methods).
188 This was also an opportunity to compare our estimates using the true trees
189 vs. the Relate-inferred trees, to examine the influence of errors in tree infer-
190 ence.

191 **Dispersal rates**

192 Our method does a good job at capturing the simulated dispersal rate when
193 using the true trees, especially at lower dispersal rates (Figure 2A). At larger
194 dispersal rates the true trees tend to cause underestimates of the simu-
195 lated dispersal rate, likely a consequence of the finite habitat we simulate
196 (with reflecting boundaries). Here we use a time cutoff of only 100 genera-
197 tions in a 50x50 habitat, meaning that with any dispersal rate larger than
198 $\sigma^2 \sim 50^2/100 = 5^2$ a lineage is reasonably likely to cross the entire habitat
199 in that time. At higher dispersal rates or longer cutoff times (Figure S1),
200 the simple Brownian motion model expects the samples to be more broadly
201 distributed than the finite habitat allows, leading to larger underestimates.
202 Regardless, we can interpret the dispersal rate inferred from the true trees
203 as a true ‘effective’ dispersal rate, given the boundaries, local competition,
204 biparental reproduction, etc. Encouragingly, the estimates from the inferred
205 trees are highly correlated with these estimates, although with an upward
206 bias. This upward bias is expected given the combination of isolation-by-
207 distance and errors in inferring tree topologies, causing geographically distant
208 samples to be mistakenly inferred to be close relatives. The bias (and vari-
209 ance across replicates) increases when we just use the panmictic coalescent
210 prior from Relate (Figure S2), showing that our spatial prior implemented
211 via importance sampling improves our inference.

212 In natural populations dispersal rates likely vary through time and so
213 we also simulated a two-epoch scenario, where the dispersal rate switched
214 100 generations ago. Figure 2B,C shows the cases where the dispersal rate
215 switched from $\sigma^2 = 0.25^2$ to $\sigma^2 = 0.5^2$, and vice-versa. As expected, our

216 estimates of the single dispersal rate fell in between the two different rates of
 217 dispersal (grey dots), averaging over the two epochs. We then applied our ex-
 218 tension of the likelihood-based inference to allow for different dispersal rates
 219 in different epochs (with *a priori* switch times between epochs, see [Materi-](#)
 220 [als and Methods](#)) to these simulations. Our multi-epoch dispersal estimates
 221 using both the true and inferred trees captured the switch in dispersal rates
 222 (when supplied with the correct switch time), again with an upward bias
 223 when using the inferred trees. Ideally we would use likelihood ratio tests to
 224 test for significantly different dispersal rates between epochs, however, our
 225 simulations show that our dispersal estimates are not well calibrated under
 226 the null model of no change (Figure S3A,D, [Supplementary Text](#)). Instead

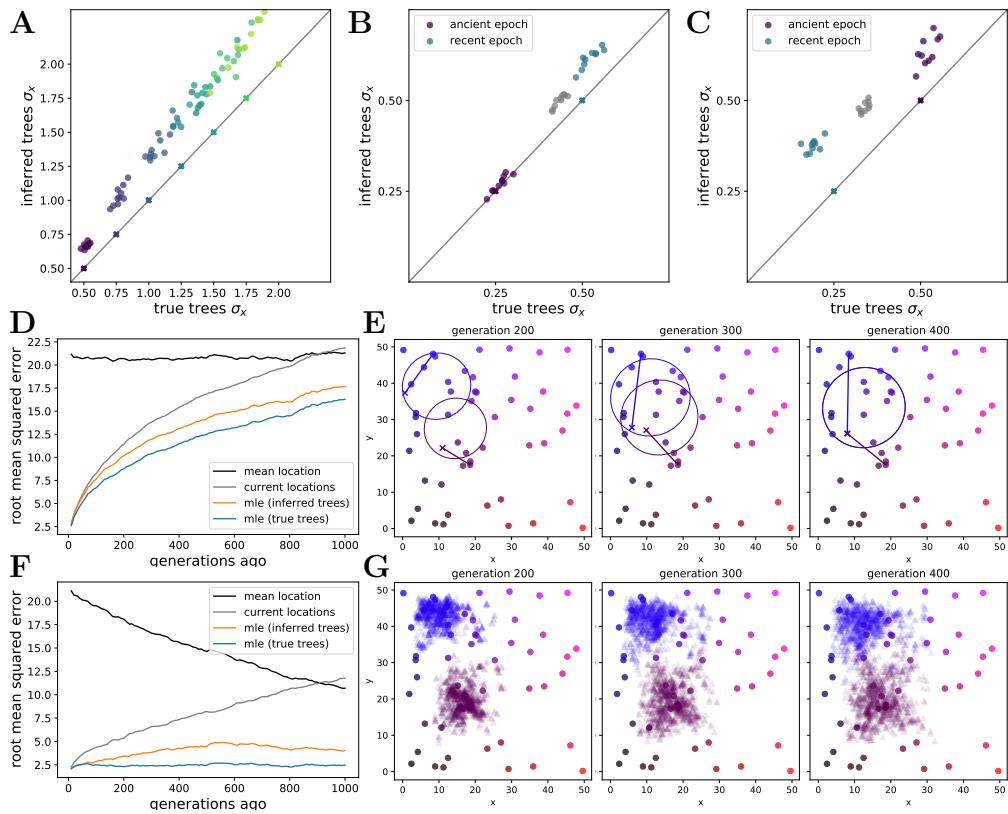


Figure 2: (Continued on the following page.)

Figure 2: **Simulations.** **(A) Accuracy of genome-wide dispersal rates.** Maximum composite likelihood estimates (over 10^2 evenly-spaced loci (trees)) of dispersal rate (only showing the standard deviation in the x dimension) using true trees vs. **Relate**-inferred trees (10 importance samples/locus, time cutoff of 10^2 generations). Each dot is a single simulation and the colours represent the simulated dispersal rates (value indicated by ‘x’ along diagonal; 10 replicates of each). **(B-C) Ability to detect time-varying dispersal rates.** Maximum composite likelihood estimates of dispersal rate using true trees vs. **Relate**-inferred trees in a two-epoch model (time cutoff of 10^3 generations). In B, the simulated dispersal rate switched from $\sigma^2 = 0.25^2$ in the deeper past to 0.5^2 in the most recent 10^2 generations. In C, the dispersal rates in the two epochs are flipped. Ten replicates are shown, with one dot of each colour for each epoch. The grey dots are the dispersal estimates under the assumption of a one-epoch model. **(D) Accuracy of locating genetic ancestors at individual loci.** Root mean squared errors between the true locations of ancestors (for 10^2 samples at 10^2 loci at 10^2 times) and the mean location of the samples, the current locations of the samples, and the maximum likelihood estimates from the inferred and true trees. **(E) Locating genetic ancestors at a particular locus.** 95% confidence ellipses for the locations of genetic ancestors for two samples at a single locus (using the true trees and the maximum likelihood dispersal rate). The ‘x’s are the true ancestral locations and the lines connect the true ancestral locations with the location of the contemporary descendant sample. **(F) Accuracy of mean genetic ancestor locations.** Root mean squared errors between the true mean (over 10^2 loci) location of ancestors (for 10^2 samples at 10^2 times) and the mean location of the samples, the current locations of the samples, and the mean (over 10^2 loci) maximum likelihood estimates from the inferred and true trees. **(G) Locating genetic ancestors at many loci.** Maximum likelihood estimates for the locations of genetic ancestors for the same two samples at 10^3 loci, here using true trees and the maximum likelihood dispersal rate. Panels D-G are simulated with $\sigma_x = \sigma_y = 0.5$ and $\rho = 0$.

²²⁷ we use the point estimates as a useful heuristic to detect broad patterns of
²²⁸ dispersal rate change.

229 **Locating ancestors**

230 We next wanted to test our ability to locate the genetic ancestor of a sam-
231 pled genome at a given locus and a given time. Our likelihood-based method
232 gives both point estimates (maximum likelihood estimate, MLE) and 95%
233 confidence ellipses (under the Brownian motion model), constructed based
234 on the MLE of the genome-wide dispersal rate. We also have developed a
235 best linear unbiased predictor (BLUP) of ancestral locations – importance
236 sampling over analytically calculated MLE locations, rather than numeri-
237 cally finding the maximum of importance sampled likelihoods – that is faster
238 to calculate, makes fewer assumptions, and is less reliant on the estimated
239 dispersal rates (and completely independent of them when only one epoch),
240 but this method does not give measures of uncertainty. Here we focus on
241 the MLE-based method for estimating ancestral locations, but both meth-
242 ods are implemented in our software and both give essentially identical point
243 estimates for ancestral locations in our simulations.

244 Figure 2D shows the error in the MLE ancestor locations, using the true or
245 inferred trees, and compares this to sensible straw-man estimates (the current
246 location of each sample and the mean location across samples). We see that
247 the true trees give the best estimates and that the inferred trees do only
248 slightly worse. That said, the mean squared error in our inferred location of
249 the ancestor at a locus grows relatively rapidly back in time. To illustrate the
250 cause of this increase in error, Figure 2D shows the 95% confidence ellipses for
251 the locations of the ancestors of two samples at one particular locus at three
252 different times in the past. In this example the lineages coalesce between
253 generation 300 and 400, so the ellipses merge. While the ellipses do a good
254 job of capturing true ancestral locations (the ‘x’s), the size of an ellipse at
255 any one locus grows as we move back in time, meaning at deeper times (and
256 higher dispersal rates) any one locus contains little information about an
257 ancestor’s location. Given the large uncertainty of an ancestor’s location at
258 any one locus, we combine information across loci and consider a cloud of
259 MLE ancestor locations from loci across the genome for a particular sample
260 at a particular time in the past (Figure 2G). The genome-wide mean of the
261 MLE locations is able to predict the true mean location of genetic ancestors
262 with much lower error (Figure 2F), even with the inferred trees, suggesting
263 our method can successfully trace major geographic ancestries of a sample
264 back into the past.

265 **Empirical application: *Arabidopsis thaliana***

266 We next applied our method to 1135 *Arabidopsis thaliana* accessions from
267 a wide geographic range ([Alonso-Blanco et al., 2016](#)). *A. thaliana* has a
268 complex, and not yet fully resolved, population history ([Fulgione and Han-](#)
269 [cock, 2018; Hsu et al., 2019](#)). The samples in this 1001 Genomes dataset are
270 thought to have descended from individuals from at least two glacial refugia,
271 including one refuge in north Africa that contributed substantial ancestry
272 to the ‘Iberian Relict’ samples ([Alonso-Blanco et al., 2016; Durvasula et al.,](#)
273 [2017; Fulgione et al., 2018](#)) and one refuge near the Balkans that contributed
274 substantial ancestry to the more weedy and cosmopolitan ‘Non-Relict’ sam-
275 ples ([Lee et al., 2017](#)). It is thought that populations from both refuges
276 first expanded northwards, followed by a fast east and west expansion of
277 the Balkan population across most of Eurasia ([Alonso-Blanco et al., 2016;](#)
278 [Lee et al., 2017; Fulgione and Hancock, 2018; Hsu et al., 2019](#)). Extensive
279 admixture between the expanding populations appears to have obscured the
280 timing of the most recent east-west expansion, with some estimates before
281 ([Durvasula et al., 2017; Fulgione et al., 2018](#)) and some after ([Alonso-Blanco](#)
282 [et al., 2016; Lee et al., 2017; Hsu et al., 2019](#)) the last glaciation, which
283 ended ~11,000 years ago. The 1001 Genomes dataset contains a relatively
284 good spatial sampling of individuals with extensive Non-Relict ancestry while
285 about 2% of the samples are considered Iberian Relicts, mostly in Spain but
286 also two samples in Morocco ([Alonso-Blanco et al., 2016](#)). The dataset also
287 contains one ‘Relict’ sample from each of Cabo Verde, the Canary Islands,
288 Sicily, and Lebanon ([Alonso-Blanco et al., 2016](#)), all of which have been found
289 to contain substantial Non-Relict ancestry ([Alonso-Blanco et al., 2016; Lee](#)
290 [et al., 2017; Zeng et al., 2017](#)). The 1001 Genomes dataset does not include
291 more recent samples from Madeira ([Fulgione et al., 2018](#)), Africa ([Durvasula](#)
292 [et al., 2017](#)), and East Asia ([Zeng et al., 2017; Zou et al., 2017](#)), which likely
293 contain ancestry from additional refugia ([Fulgione and Hancock, 2018; Hsu](#)
294 [et al., 2019](#)). We did not attempt to include these more recent samples in
295 this application because, while an important part of the puzzle, their spatial
296 sampling is relatively sparse. Finally, the dataset includes 125 samples from
297 across North America, a range expansion resulting of multiple recent human
298 introductions ([Exposito-Alonso et al., 2018; Shirsekar et al., 2021](#)).

299 ***A. thaliana* tree sequences**

300 We used **Relate** (Speidel et al., 2019) to infer the tree sequence, estimate
301 the panmictic effective population size through time, and resample branch
302 lengths for importance sampling (see Materials and Methods). After drop-
303 ping the 50% of trees with the fewest number of mutations, per chromosome,
304 the tree sequence contains 213,481 trees across the 5 autosomes. The tree
305 sequences (in both anc/mut and tskit formats) and the code to generate
306 them are publicly available [link], which we hope will facilitate additional
307 analyses (e.g., inferring selection (Stern et al., 2019, 2021)).

308 *A. thaliana* is a selfer with a relatively low rate of outcrossing (Bomblies
309 et al., 2010; Platt et al., 2010), thus it is worth taking a moment to consider
310 the impact of selfing on our inferences. We chose the 1001 Genomes panel
311 because of its large size and broad geographic sampling. Further, the avail-
312 ability of inbred accessions means that the samples have been well-studied
313 and from our perspective remove the complications of obtaining phased hap-
314 lotypes to run **Relate**. On the other hand, the high rate of selfing lowers the
315 effective recombination rate and so is expected to increase the correlation in
316 genealogies along the genome (Nordborg, 2000). However, in practice, linkage
317 disequilibrium breaks down relatively rapidly in *A. thaliana*, on the scales of
318 tens of kilobases (Kim et al., 2007), such that many trees along the genome
319 should be relatively independent from each other. A related issue is that
320 the individuals with recent inbreeding (selfing) in their family tree will have
321 fewer genealogical and genetic ancestors than outbred individuals. Thus in
322 any recent time-slice there are a reduced number of independent genetic an-
323 cestors of a individual from a selfing population, but even with relatively low
324 rates of outbreeding the number of ancestors still grows rapidly (Lachance,
325 2009). Finally, while the effective recombination rate may vary through time
326 along with rates of selfing, **Relate** uses a mutational clock to estimate branch
327 lengths, and thus they should be well calibrated to a generational time scale.

328 **Rapid recent east-west dispersal**

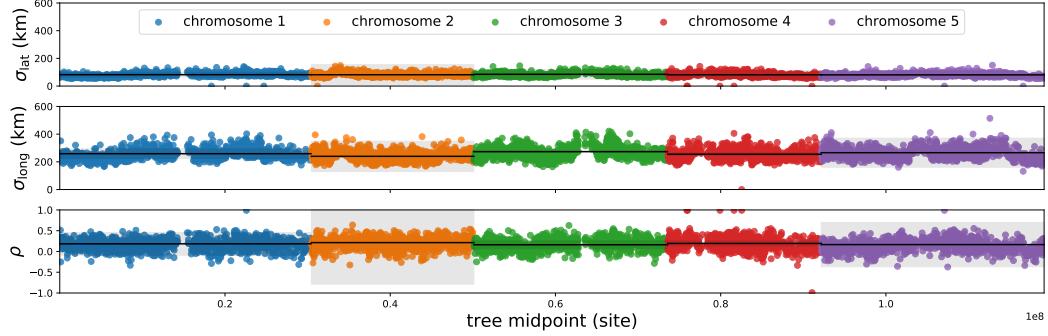
329 We first used the tree sequence to estimate dispersal rates. In doing this
330 we ignored the locations of 199 genetically near-identical samples ($< 10^3$
331 basepairs different, as in Alonso-Blanco et al., 2016), the 125 samples from
332 North America, and the 2 samples from Japan (see Materials and Methods).
333 The latter two groups are outliers for dispersal as they are thought to have

334 been carried large distances by humans in the recent past ([Exposito-Alonso](#)
335 [et al., 2018; Zou et al., 2017](#)). The near-identical samples may also rep-
336 resent recent rare long-distance dispersal, but could alternatively be due to
337 mis-assignments or mix-ups after collection ([Alonso-Blanco et al., 2016](#)), and
338 so we remove them in case of the latter. Removing near-identical samples
339 that are true long-distance migrants will cause us to underestimate dispersal
340 rates. However, most near-identical samples are very near one another and
341 are likely from the same inbred lineage ([Alonso-Blanco et al., 2016](#)); exclud-
342 ing these will have little effect on our inference. In estimating dispersal we
343 use 10 importance samples of branch lengths per tree and a time cutoff of
344 10^4 generations (equivalently, years). Estimates with a cutoff of 10^3 were
345 essentially identical, suggesting most of the information on dispersal comes
346 from movements in the last 10^3 years, but having trees that go back 10^4
347 generations allows us to locate ancestors more reliably at deeper times.

348 We first estimate a single, constant dispersal rate, i.e., assuming one
349 epoch. Figure 3A shows both the per-locus (dots) and the composite per-
350 chromosome (horizontal lines, with gray denoting the standard error) es-
351 tates. We find that the per-generation rate of dispersal is ~ 10 times
352 larger across longitude than across latitude (i.e., $\sigma_{\text{long}}^2 \sim 10\sigma_{\text{lat}}^2$, with units
353 $\text{km}^2/\text{generation}$). This high rate of longitudinal dispersal is consistent with
354 the hypothesis of rapid expansion along the east-west axis of Eurasia from
355 glacial refugia, facilitated by relatively weak environmental gradients and,
356 potentially, human movements and disturbance ([Alonso-Blanco et al., 2016](#);
357 [Lee et al., 2017](#)). As discussed in [Materials and Methods](#), our approach does
358 not condition on the sample locations, which means that the distribution
359 of sample locations may influence dispersal estimates. This is a particular
360 worry here where we have samples from a wider range of longitudes than lat-
361 itudes and we find a larger longitudinal dispersal rate. Fortunately, we can
362 check our result by estimating dispersal separately in longitude and latitude
363 (assuming no covariance), where conditioning on sample locations is much
364 weaker ([Meligkotsidou and Fearnhead, 2007](#)). Doing this we find essentially
365 identical dispersal rates as before (mean dispersal rates across chromosomes:
366 $\sigma_{\text{long}} \approx 259 \text{ km}/\sqrt{\text{gen}}$, $\sigma_{\text{lat}} \approx 81 \text{ km}/\sqrt{\text{gen}}$), supporting our finding of faster
367 dispersal across longitude ($\sigma_{\text{long}}^2 \sim 10\sigma_{\text{lat}}^2$).

368 To explore the idea that environmental changes, such as human move-
369 ments and disturbance, have affected the rate of spread of *A. thaliana*, we
370 next estimated per-chromosome dispersal rates under a two-epoch model.

A) One-epoch model (per-locus and per-chromosome estimates)



B) Multi-epoch models (per-chromosome estimates only)

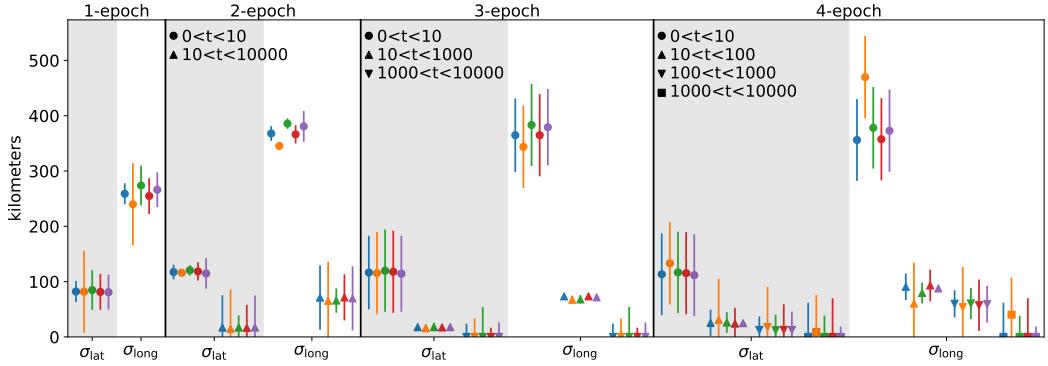


Figure 3: Dispersal rate estimates in *Arabidopsis thaliana*. **(A) One-epoch model.** The dots are maximum likelihood estimates of dispersal rate (σ , in units of $\text{km}/\sqrt{\text{generation}}$) at 10^3 evenly-spaced loci (trees) per chromosome. The black lines are maximum composite likelihood estimates for each chromosome (using 10^2 evenly-spaced loci) and the grey shading is the standard error (estimated from the Hessian at the maximum). **(B) Multi-epoch models.** The maximum composite likelihood dispersal rate estimates (\pm standard error) across each chromosome (10^2 evenly-spaced loci) for models with multiple epochs (and the one-epoch model again for comparison). For two- and three-epoch models with alternative (and less likely) split times, see Figures S4 and S5.

371 Regardless of whether the more recent epoch ends 10, 10^2 , or 10^3 generations
 372 ago, dispersal rates in the recent epoch are greater than those under the one-

epoch model (Figure S4). Figure 3B shows the dispersal estimates under the split time with the highest likelihood (10 generations ago), suggesting very rapid dispersal along the east-west axis in very recent times (on the order of decades). We also ran three-epoch models that allow dispersal to change at 10 and 10^2 , 10 and 10^3 , or 10^2 and 10^3 generations ago (Figure S5). The model with the highest likelihood (split times of 10 and 10^3 generations ago, Figure 3B) further supports the idea that the signal of very rapid dispersal comes from movements on the timescale of decades. Finally, we ran a four-epoch model with split times of 10, 10^2 , and 10^3 generations ago (Figure 3B). This had the highest likelihood of any model we ran and reiterates the main conclusion – there appears to have been very rapid recent east-west dispersal. While it is tempting to compare dispersal rates across epochs, one caveat here is that finite habitat boundaries may disproportionately depress more ancient estimates; in essence, while we need a large recent dispersal rate to account for the large distances between relatively closely related samples, if such a dispersal rate continued far into the past the samples would be expected to cover a much wider geographic range than is possible given the constraints of habitat.

Identifying interesting dispersal outliers

We next used the tree sequence and dispersal estimates to locate genetic ancestors (see Materials and Methods). We can locate ancestors at every locus for every individual at any time, which represents an incredibly rich resource for understanding population movement. As a first step, we visualize the mean ancestral location, averaged over loci, for every individual to detect samples with unusual geographic ancestries. To do this we estimated the locations of recent ancestors of all samples at 100 evenly-spaced loci per chromosome and averaged over loci to give a mean ancestral ‘displacement’ (from the sample, backwards in time) for each sample.

Figure 4 shows these mean displacements as arrows, with colours emphasizing the length of the arrow, at 10 and 100 generations ago. As expected, most arrows point inwards, corresponding to the ancestors of the sample being geographically closer to one another than the samples are. There are a few exceptions however. For example, there is a sample in Romania (accession 9737) with a mean displacement far to the east. This outlier appears to coalesce at many loci ~ 100 generations ago with two samples in Russia near Kazakhstan’s northeastern border (accessions 9627 and 9630), reflecting a

409 recent long-distance dispersal event. Taking a look at the coalescence times
410 between accessions 9737 and 9627 along the tree sequence (Figure S6A) there
411 are many large blocks of recent coalescence (< 100 generations), covering
412 about 50% of the genome, while the remainder of the genome coalesces more
413 deeply. This is consistent with the Romanian sample's previous placement in
414 the 'Asia' admixture group (Alonso-Blanco et al., 2016). A second example
415 is a sample from southern France (accession 9933) that quickly moves east to
416 Romania/Ukraine. Taking a look at coalescence times between this sample
417 and those further to the east, we find this sample coalesces < 100 generations
418 ago with a sample from Afghanistan (accession 10015) over a few very large
419 blocks of its genome, but is much further diverged elsewhere (Figure S6B).
420 This is consistent with this French sample's previous 'Admixed' admixture
421 assignment (including substantial ancestry from the 'Central Asia' group;
422 Alonso-Blanco et al., 2016).

423 There are also some samples that are outliers in terms of distance traveled
424 by their ancestors. For example, there are two samples, one from the UK
425 (accession 7314) and the other from Belarus (accession 6981), that travel long
426 distances towards one another in the past \sim 10 generations (see also Figure
427 3B in (Alonso-Blanco et al., 2016)). A look at the coalescence times between
428 these two samples shows that these two individuals are extremely closely
429 related along almost all of their genome except the first \approx 700kb of chro-
430 mosome 5 (Figure S6C). Other outliers include the sample from Cabo Verde
431 (accession 6911), which is classified as a Relict but includes a large amount
432 of European ancestry (Alonso-Blanco et al., 2016) and the most eastward
433 Russian sample (accession 9622), which has substantial European ancestry
434 (Alonso-Blanco et al., 2016). More generally, we also see greater rates of
435 movement, on average, from samples from further east, consistent with the
436 previous finding of reduced genetic distance between 'Asian' samples up to
437 \sim 250km apart (Alonso-Blanco et al., 2016), the signal of a rapid eastward
438 expansion.

439 **Detecting the existence and source of rare long-distant migrants
440 and misplaced samples**

441 We next looked in more detail at some of the largest outliers in recent an-
442 cestral movements. As mentioned above, the two accessions from Japan are
443 thought to be recent long-distance migrants. When we include these sam-
444 ples' locations and locate their ancestors, these samples are the largest out-

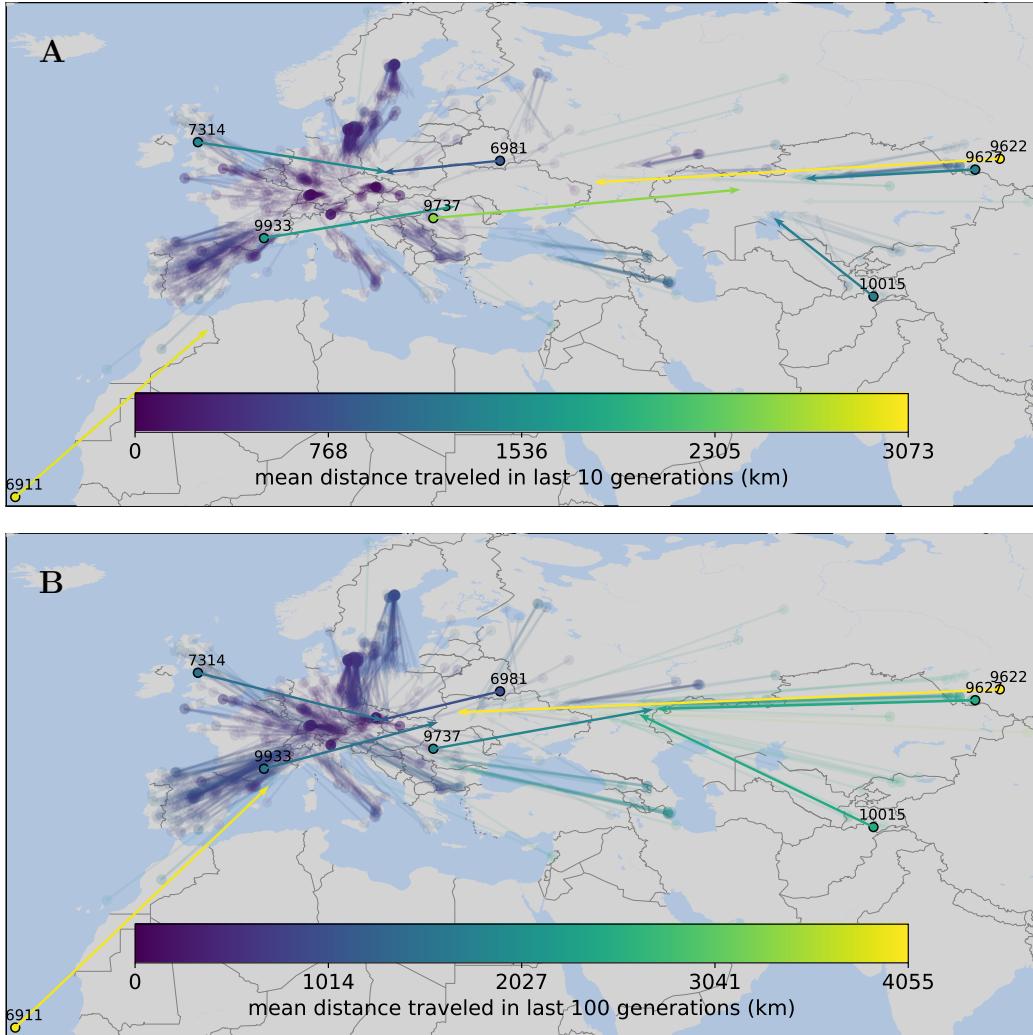


Figure 4: **Locating genetic ancestors to identify interesting outliers.** Mean maximum likelihood genetic ancestor locations (averaged over 10^2 evenly-spaced loci per chromosome) (A) 10 and (B) 100 generations ago. The samples discussed in the text are highlighted and labeled (accession numbers). Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

445 liers in their inferred rate of recent movement (result not shown). Creating a
446 smooth kernel from 10^3 inferred ancestral locations per chromosome, Figure
447 5A shows an extremely rapid transition between Japan and Europe in the
448 last $\sim 10 - 100$ generations (only accession 7207 shown). This is consistent
449 with the hypothesis that humans moved the ancestors of these samples from
450 Europe to Japan in the last few hundred years (Zou et al., 2017).

451 The Japanese samples, while spatial outliers, are relatively closely related
452 to the majority of the other samples (Alonso-Blanco et al., 2016). Thus, our
453 method has the power to locate their ancestors and detect them as recent
454 migrants. To see if the our method, and the data, have the power to de-
455 tect a recent migrant from a less well represented ancestry, we deliberately
456 misplaced the location of an Iberian Relict from Spain (accession 9533) to
457 the mean location of the Eurasian samples after filtering (southern Czech
458 Republic). Since this is an Iberian Relict, it is most closely related to the 21
459 other Iberian Relict samples (only $\approx 2\%$ of the samples), with much of its
460 ancestry likely coming from a refuge in north Africa (Alonso-Blanco et al.,
461 2016; Durvasula et al., 2017; Fulgione et al., 2018). Figure 5B shows that,
462 despite the rarity of this ancestry in this dataset, our method correctly puts
463 the misplaced Relict back into Spain at most loci in ~ 10 generations.

464 The examples above show that our method has the power to detect
465 many recent long-distant migrants and misplaced samples, and identify their
466 source. Note that, because we assume a model of continuous migration, the
467 ancestors of recent migrants and misplaced samples have to migrate through
468 intermediate locations to reach their likely source location. However, if, after
469 inspection, an investigator suspects that a sample was likely misplaced, or
470 dispersed a long distance suddenly, they could choose to ignore that sample's
471 location and use the trees alone to estimate its source.

472 Inferring the origin and ancestry of a recent invasion

473 To illustrate how we can use our method to identify the source location of
474 a subset of samples (e.g., recent migrants or misplaced samples), we used
475 our method to infer the origin and ancestry of the 125 accessions from the
476 recent expansion of *A. thaliana* into North America. We did this by ignor-
477 ing the locations of these samples and predicting their current and ancestral
478 locations based only on the trees and the locations of the remaining sam-
479 ples. Figure 6 shows the predicted current (purple) and ancestral locations,
480 averaged over 10^2 loci per chromosome, for each North American accession.

481 This does two things. First, if we knew the time of colonization, we could
482 read off where we expect the colonization to have originated from. Instead
483 of simply asking where are the most closely related samples, we allow all the
484 lineages to move backwards in time, naturally correcting for the movement

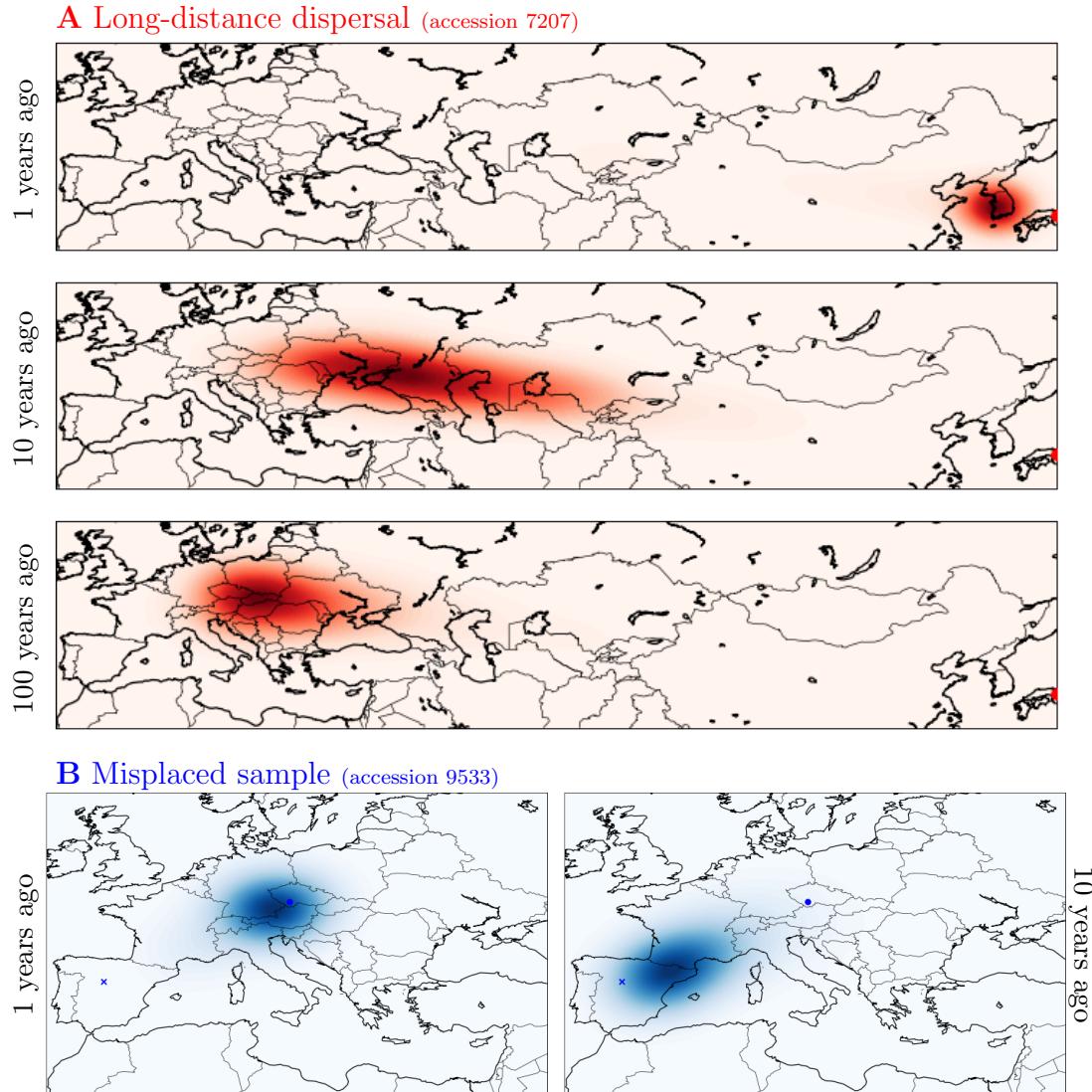


Figure 5: (Continued on the following page.)

Figure 5: **Locating genetic ancestors to detect the existence and source of (A) recent long-distance dispersal and (B) misplaced samples.** Shown are density kernels of the maximum likelihood ancestral locations at 10^3 evenly-spaced loci per chromosome for (A) a Japanese accession (100% ‘Central Europe’ ancestry) and (B) an Iberian Relict (100% Relict ancestry) from Spain that we misplaced to the mean location of all samples (Czech Republic). The dots show the assumed sample locations (at the edge of the panes in A) and the ‘x’ in B shows the true sample location. Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

of lineages between the time of sampling and the time of colonization. In this case, the colonization of North America by *A. thaliana* is thought to be about 400 years ago ([Exposito-Alonso et al., 2018](#)), suggesting that the majority of the North American accessions came from southern Germany, despite being more closely related, in many cases, to samples in north eastern France today. Second, Figure 6 helps visualize the high degree of relatedness among the North American samples. For instance, a large group of closely related North American samples are also closely related to current-day samples in north eastern France, and appear to coalesce with one another at many loci in the last few hundred years, near the purported time of colonization. On the other hand, there is a handful of more distantly related North American samples, most closely related to present-day samples ranging from northern Germany to Russia, that do not coalesce at most loci in the past 10^3 years, suggesting multiple colonizations of North America, as previously hypothesized ([Exposito-Alonso et al., 2018](#)). Encouragingly, connecting the sample and inferred locations (Figure S7), we find that the samples placed in Russia and Slovakia are from near Manistee, Michigan (accessions 1890 and 2202) and the sample placed near northern Germany is the reference, Col-0, from Missouri (accession 6909), consistent with recent independent findings ([Shirsekar et al., 2021](#)).

505 Visualizing alternative geographic ancestries

As mentioned above, the 1001 Genomes dataset is thought to contain samples drawing ancestors from at least two different glacial refugia, including the Non-Relict samples with ancestry predominately from a Balkan refuge ([Lee](#)

509 et al., 2017) and the Iberian Relict samples with substantial ancestry from
510 a refuge in north Africa (Alonso-Blanco et al., 2016; Durvasula et al., 2017;

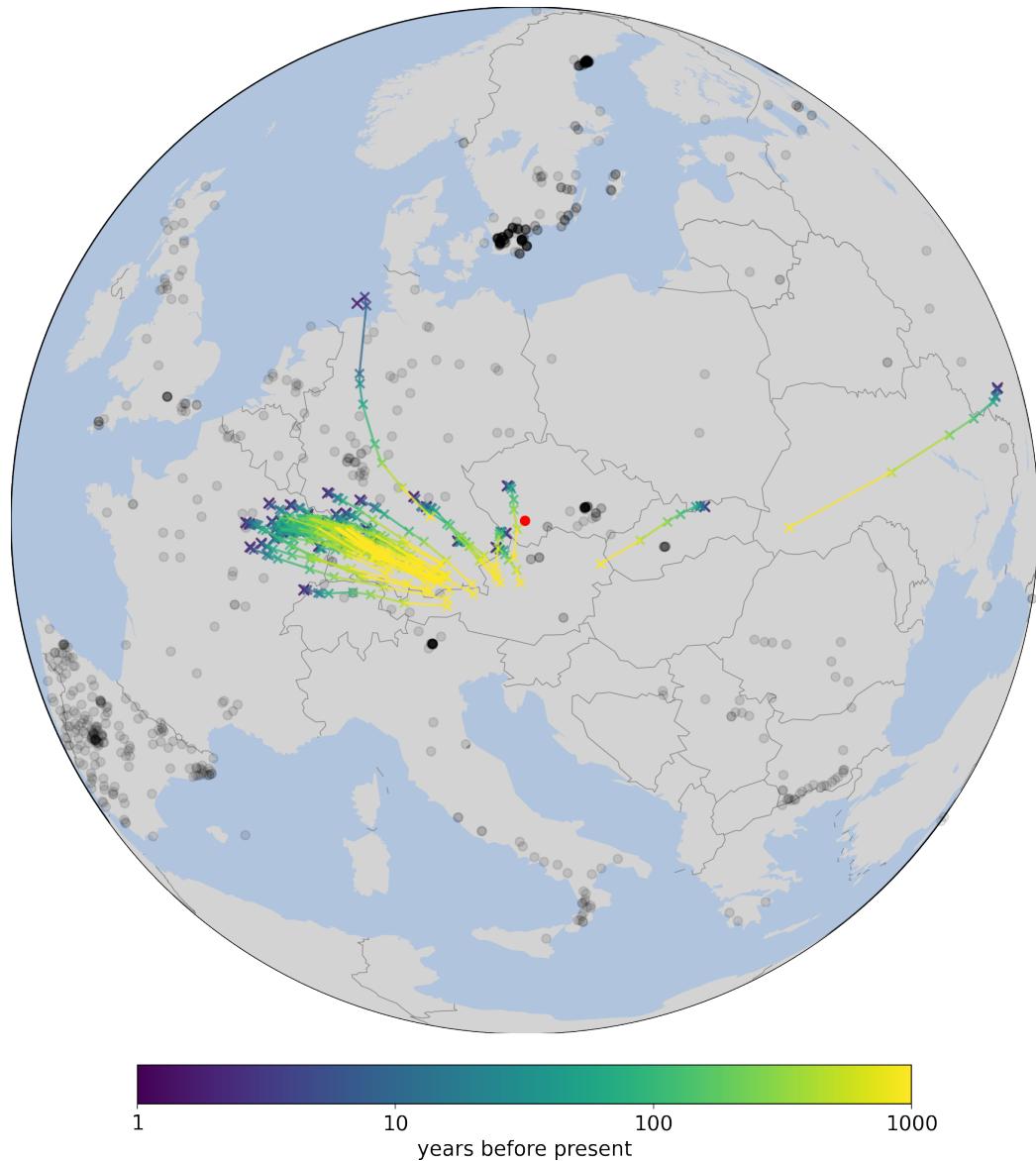


Figure 6: (Continued on the following page.)

Figure 6: **Locating genetic ancestors to identify the source and ancestry of recent expansions/colonizations.** Mean maximum likelihood genetic ancestor locations (averaged over 10^2 evenly-spaced loci per chromosome) of the 125 North American accessions over the last 10^3 years. The black dots are sample locations and the red dot is their mean. Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

511 [Fulgione et al., 2018](#)). We therefore sought to use our method to visualize
512 these alternative geographic ancestries.

513 We first compared the geographic ancestries of three representative ac-
514 cessions from Spain (Figure 7): one from a Non-Relict admixture group (ac-
515 cession 6933; 100% ‘Spain’ ancestry); one from the Relict admixture group
516 (accession 9533; 100% Relict ancestry); and one drawing roughly equal an-
517 cestry from both these groups (accession 9530; 57% Relict and 43% ‘Spain’
518 ancestry). As expected given the glacial refuge of the Non-Relict lineages
519 is thought to be near the Balkans ([Lee et al., 2017](#)), the Non-Relict sam-
520 ple’s genetic ancestors move gradually and coherently north east, out of the
521 Iberian Peninsula by 200 years ago and into northern Italy/Austria by 500
522 years ago. In contrast, a large proportion of the Relict sample’s ancestors
523 remain on the Iberian Peninsula for the last 500 years, where the rest of the
524 closely related Iberian Relict samples are clustered. The Admixed sample’s
525 genetic ancestors display much more variability, with the rate of return to
526 the east depending on the ancestry at each locus, with many ancestors still
527 remaining in Spain but many already in Italy \approx 200 years ago.

528 To confirm this is a general pattern beyond these three representative
529 samples, we also looked at the mean ancestral displacements (over 10^2 loci
530 per chromosome) of all mainland Spanish samples from each of the three
531 groups (Relict, Non-Relict, and Admixed; Figure 8). We see that, on average,
532 Non-Relict samples tend to move both further north and further east than
533 the Relict samples, in the direction of the mean sample location and towards
534 the Balkans (the same conclusions are reached if we compare just the ‘Spain’
535 and Relict groups, results not shown). The Admixed samples show much
536 more variation in the amount of northward displacement, as expected given
537 these samples could have ancestries from any two or more groups, including
538 two Non-Relict groups (e.g., ‘Spain’ and ‘North Sweden’).

539 Eventually nearly all of an Iberian Relict sample’s ancestors are pulled

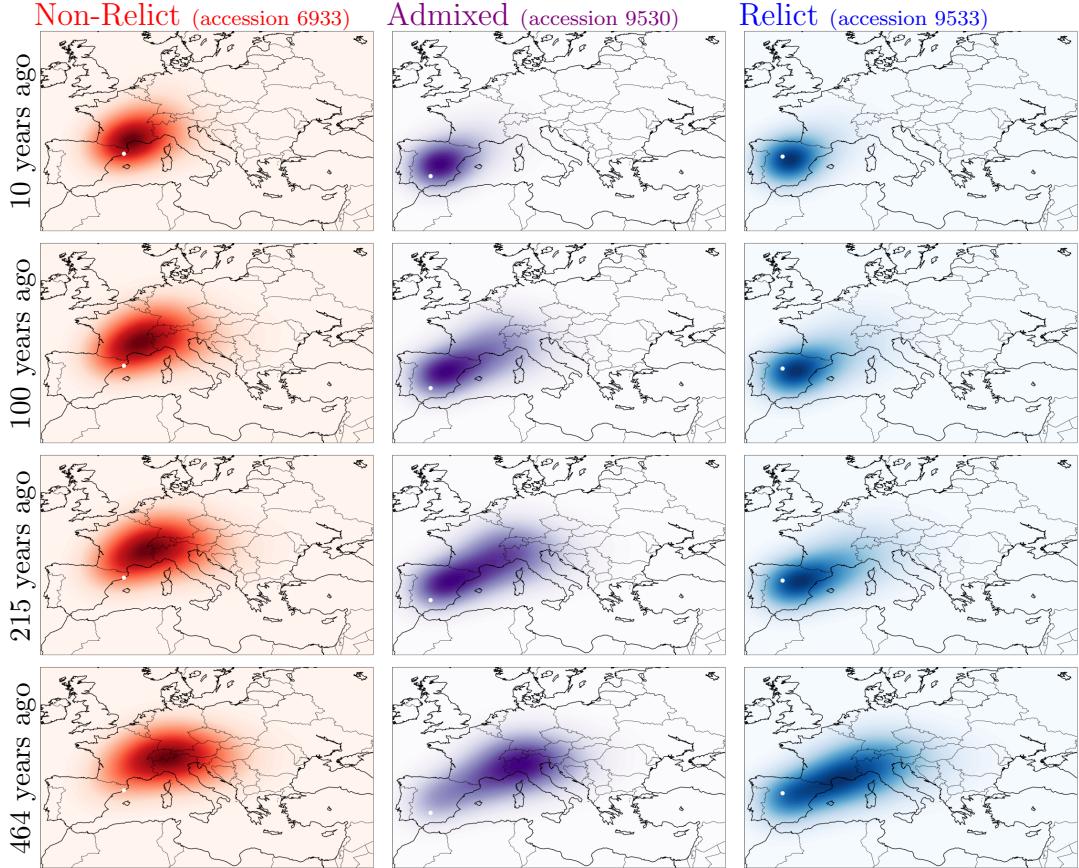


Figure 7: Locating genetic ancestors to depict alternative geographic histories. Genetic ancestors of three samples from Spain from three different admixture groups: a Non-Relict (100% ‘Spain’), an Admixed (43% ‘Spain’, 57% Relict), and an Iberian Relict (100% Relict). Shown are density kernels of the maximum likelihood ancestral locations using 10^3 evenly-spaced loci per chromosome. White dots show the sample locations. Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

east, perhaps because the Iberian Relict ancestry is only represented by $\sim 2\%$ of the samples, but perhaps also because the Iberian Relicts are through to have extensive admixture with the Non-Relicts (Fulgione et al., 2018). There

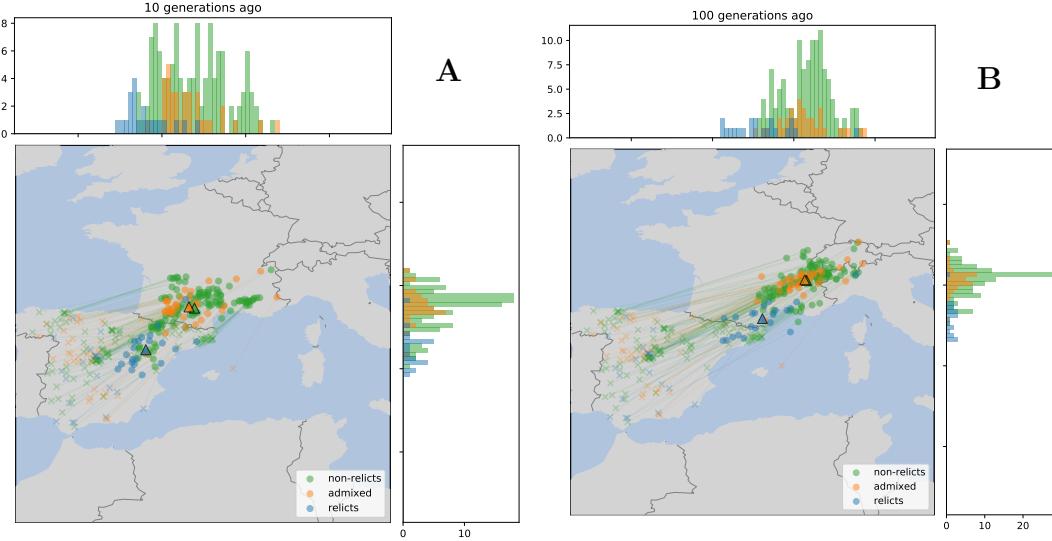


Figure 8: Relict samples display an alternative geographic history. Mean genetic ancestor locations of all samples from Spain (excluding the Canary Islands) by admixture group. Data as in Figure 4, with ‘x’s the sample locations and ‘o’s the mean ancestor location and lines connecting them. Colours indicate admixture group (Non-Relict is any group besides Relict or Admixed; Admixed is any sample with < 60% ancestry in all groups ([Alonso-Blanco et al., 2016](#))). Triangles indicate mean locations for each admixture group. Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

543 are only three samples from Africa in the 1001 Genomes dataset (accessions
 544 6911, 9606, and 9939), all of which show admixture with European groups
 545 and are inferred to have geographic ancestries eventually tracing north east
 546 (Figure 4). It would be interesting to see if adding the ~ 80 more African
 547 genomes currently available ([Durvasula et al., 2017](#)) would pull some of the
 548 ancestors of these Iberian Relicts towards a possible north African refuge
 549 ([Durvasula et al., 2017; Fulgione et al., 2018](#)), and help better visualize the
 550 Relict-Non-Relict admixture.

551 **Visualizing the sources of admixture**

552 As we have seen, the 1001 Genomes dataset contains a number previously
553 identified Admixed samples (Alonso-Blanco et al., 2016). This provides a
554 nice opportunity to explore how locating recent ancestors can help visual-
555 ize admixed ancestry and its geographic sources. To do this we plot recent
556 ancestral displacements, as in Figure 4, but this time without averaging
557 over loci, instead plotting the displacement for each locus. Figure 9A shows
558 the inferred displacements (in the past 10 generations) for 6 of the more
559 striking Admixed samples, both in terms of the direction and magnitude
560 of ancestral displacements. To help summarize the range of displacements
561 from each sample, the ‘windrose’ insets show a histogram of displacement
562 directions, weighted by displacement length and coloured by chromosome.
563 In many cases there is a large spread in the direction of the displacements,
564 illustrating multiple contributing geographic ancestries. For example, the
565 sample from France discussed earlier (accession 9933) appears to coalesce
566 recently with lineages from both relatively near north east and relatively far
567 east (consistent with its previous admixture assignment containing substancial
568 ‘Germany’ and ‘Central Asia’ ancestry) and the sample from Romania
569 (accession 9743) appears to get ancestry from both the west and east (con-
570 sistent with its previous admixture assignment containing substancial ‘Italy,
571 Balkans, & Caucasus’ and ‘Central Asia’ ancestry).

572 There is also information contained in the correlation of ancestral move-
573 ments across the genome. To demonstrate how geographic ancestries are
574 correlated across the genome, Figure 9B shows each displacement emanating
575 from it’s location along a linear representation of the genome (rather than all
576 emanating from the same place, as in the map) for the French and Romanian
577 sample (see Figure S8 for the other 4 samples). With this visual we see that
578 the French sample gets much of its most eastern ancestry from chromosomes
579 3 and 5 (green and purple) while the Romanian sample gets most of its east-
580 ern ancestry from chromosome 4 (red). This figure also serves to illustrate
581 just how much information is being inferred – at any time in the past, we
582 can estimate the locations of thousands of genetic ancestors for thousands
583 of samples, providing a rich source of information to explore and use to test
584 hypotheses.

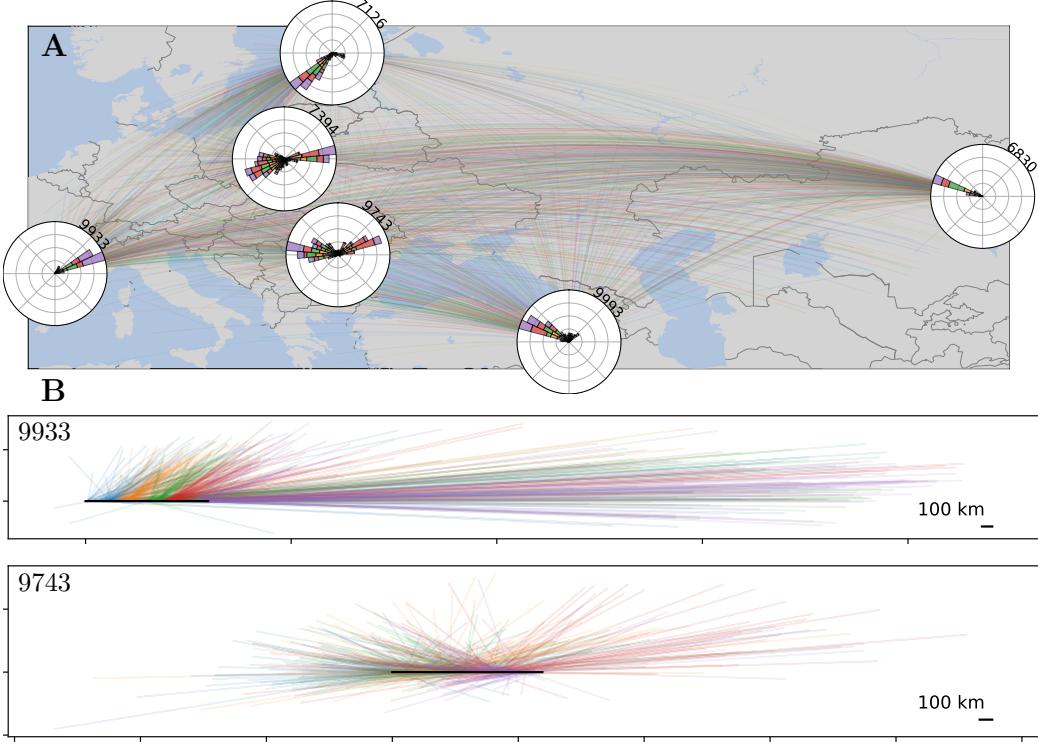


Figure 9: Locating genetic ancestors to visualize admixture and its sources. (A) The curves on the map are great circle paths connecting the sample locations (center of ‘windrose’ insets) to the maximum likelihood ancestor locations 10 generations ago at 10^2 evenly-spaced loci per chromosome for 6 striking Admixed samples (accession numbers as labels). The windrose insets show a histogram of directions from the sample to the ancestors, weighted by distance. Colours indicate chromosomes, as in Figure 3. (B) The black horizontal line represents the genome and the coloured lines show the ancestral displacements in longitude and latitude (converted to kms) from that position of the genome. See Figure S8 for the genome-view of ancestral displacements of all 6 samples. Maximum likelihood locations calculated using per-chromosome four-epoch dispersal rates (Figure 3B).

585 Discussion

586 Summary of main results

587 We have developed a method that uses a sequence of trees along a recombin-
 588 ing genome – a genome-wide genealogy²⁷ – to estimate individual-level disper-

sal rates and locate genetic ancestors (Figure 1). At the core of our method is a simple model of Brownian motion, allowing likelihoods to be quickly calculated from shared evolutionary times and sample locations. On top of this we layer on importance sampling to correct for bias in inferred branch lengths, a time cutoff to ignore strong violations of the model in the deep past, and multiple epochs to allow the dispersal rate to vary through time. Simulation tests show that our method can accurately infer dispersal rate (with a slight overestimate caused by errors in tree inference), detect shifts in the dispersal rate over time, and trace major geographic ancestries hundreds of generations into the past (Figure 2). Applying our method to more than one thousand *Arabidopsis thaliana* genomes across a huge geographic expanse (Alonso-Blanco et al., 2016), we find strong evidence for very rapid recent dispersal (Figure 3), especially across longitude, and show how locating genetic ancestors can detect and visualize *i*) long-distance dispersal events (Figures 4–5), *ii*) the source of population expansions (Figure 6), *iii*) alternative geographic ancestries (Figures 7–8), and *iv*) admixture (Figure 9).

Comparison with previous methods

The idea of using trees to estimate continuous ancestral characters and their rate of change is an old one. This was originally applied to population-level characters, such as the frequency of genes in a population and their rate of genetic drift, in a phylogenetic context (Cavalli-Sforza et al., 1964; Cavalli-Sforza and Edwards, 1967; Edwards, 1970; Felsenstein, 1973, 1985; Grafen, 1989). DNA sequencing later allowed the inference of a single tree relating individual samples for a sufficiently long non-recombining sequence (as found on the human Y chromosome, in a mitochondrial genome, or in the nuclear genome of a predominately non-recombining species), which led to estimates of dispersal rates and ancestral locations under the banner of ‘phylogeography’ (Avise, 2009; Knowles, 2009). Phylogeography has proven incredibly useful, especially to infer the geographic origin and spread of viruses (Biek et al., 2007; Lemey et al., 2009, 2010; Bedford et al., 2010; Volz et al., 2013), such as SARS-CoV-2 (Worobey et al., 2020; Lemey et al., 2020; Dellicour et al., 2021a,b; Martin et al., 2021). Extending phylogeography to frequently-recombining sequences is not straightforward as there are then many true trees that relate the samples. Only recently has it become feasible to infer the sequence of trees, and their branch lengths, along a recombining

625 genome (Rasmussen et al., 2014; Speidel et al., 2019; Wohns et al., 2021).
626 Our method capitalizes on this advance to use some of the enormous amount
627 of information contained in a tree sequence of a large sample in a recombining
628 species.

629 A related method was recently demonstrated by Wohns et al. (2021),
630 who inferred the geographic location of coalescent nodes in a `tskit` tree
631 sequence (Kelleher et al., 2018), where information about nodes and branches
632 are shared between trees. Our approach differs from theirs in a number of
633 ways. First, they utilize the sharing of information about nodes and branches
634 across trees to very efficiently geographically locate every node exactly once,
635 by locating a ‘parent’ node at the midpoint between its two ‘child’ nodes
636 geographic locations and iterating up the entire tree sequence simultaneously
637 (rather than up each local tree individually). In addition to being fast this
638 has the advantage of using information from all the trees in a tree sequence.
639 In contrast, we locate ancestors independently at each local tree we consider.
640 As some ancestors (represented as nodes and branches) are shared between
641 nearby trees along the genome (though we don’t know precisely for how long
642 since we lose that information when converting the `Relate`-inferred genealogy
643 into a tree sequence), we avoid locating the same ancestors multiple times by
644 using only a sample of the trees (e.g., for *A. thaliana* we uniformly sampled
645 5000 of the ~200,000 trees, ~2.5%). Our approach therefore uses less of the
646 information in the tree sequence, in this sense. (Note that while we choose
647 trees that have low linkage disequilibrium with one another and are therefore
648 essentially unlinked, in the very recent past they will share ancestors but will
649 quickly become independent (Wakeley et al., 2012).) On the other hand,
650 when locating an ancestor we not only use information ‘below’ this ancestor
651 (i.e., its descendants’ locations and relations) but also the information ‘above’
652 the ancestor, due to the ancestor’s lineage sharing evolutionary time and a
653 recent common ancestor with non-descendant lineages. Further advantages
654 of our method include the ability to estimate dispersal rates and uncertainties
655 in ancestor locations (as we have taken a parametric approach), the ability
656 to locate ancestors at any time (not just at coalescent nodes but also at any
657 point along a branch), and accounting for uncertainty in the branch lengths
658 (using importance sampling; this could be extended to capture uncertainty
659 in the topologies as well once they can be efficiently sampled).

660 An exciting possibility is a merger of these two methods, to efficiently use
661 information from all the trees in the tree sequence and all the information
662 above and below ancestors. One approach might be to use something like

663 the inside-outside algorithm in (Wohns et al., 2021), which they use for node
664 times but not locations. A second approach would be to model Brownian
665 motion on the full ancestral recombination graph (e.g. using a program like
666 ARGweaver; Rasmussen et al., 2014). The latter approach would allow one
667 to visualize the geographic splitting and coalescence of the ancestors of a
668 genome backwards in time.

669 Future directions

670 We have chosen to stick with a very simple model of Brownian motion with a
671 piecewise-constant dispersal rate over time. There are a number of extensions
672 that could readily be applied. For instance, we could allow dispersal rates
673 to vary between branches (O’Meara et al., 2006), e.g., to compare disper-
674 sal rates in different parts of a species’ range. Or we could model dispersal
675 under a more complex model, like the Early Burst (Harmon et al., 2010)
676 or a Lévy process (Landis et al., 2013), which may help identify periods of
677 range expansion or sudden long-range dispersal. Alternatively we could take
678 a Bayesian approach, allowing much greater flexibility and the incorporation
679 of many recent advances in phylogeography. For example, one could then
680 model dispersal as a relaxed random walk (Lemey et al., 2010), which may
681 be more appropriate for sample locations that are very non-normal and could
682 incorporate habitat boundaries. Second, given that there is large variance in
683 the inferred locations of distant ancestors at any one locus (Schluter et al.,
684 1997), but very many loci, we could take an ‘empirical Bayes’ approach and
685 use the posteriors on ancestral locations over many loci to set a prior for a
686 given locus. This might be especially helpful at deeper times, e.g., tracing
687 human ancestors back 100s of thousands of years, where the noise in the
688 ancestral location at any one locus is large, yet we can be relatively certain
689 that the majority of lineages are in Africa. We might alternatively set pri-
690 ors to test hypotheses, e.g., if we surmise there were multiple glacial refugia
691 during the last glaciation we can set a prior on ancestral location with peaks
692 at these hypothesized locations and infer what percent of a sample’s lineages
693 descended from each. Models of ancestral locations based on past climatic-
694 and ecological-niche models could provide a rich source of data for building
695 such priors, and given the large amounts data available in recombining se-
696 quences these models could be subject to rigorous model choice. Finally, it
697 might also be interesting to take a more model-agnostic approach and use
698 machine learning. For example, Locator (Battey et al., 2020) uses deep

699 neural networks to infer the locations of extant individuals from unphased
700 genotype data. In essence this means **Locator** is both determining the rela-
701 tionships between samples and locating them. Separating these two steps by
702 first inferring a tree sequence and then supplying information from this tree
703 sequence to the deep neural network may both improve location estimates
704 for extant individuals and allow such a method to locate genetic ancestors.

705 Our approach relies on the locations of the current-day samples. While
706 we have shown that we can learn much about the geographical history of a
707 species with this approach, its accuracy is necessarily limited. For example,
708 if historical parts of the range are undersampled in a particular dataset the
709 method will struggle to locate ancestors in these regions, particularly further
710 into the past, as we saw with the Iberian Relict samples with ancestry from
711 the putative north African refugia. Similarly, if a species' range has shifted
712 such that few present day individuals exist in portions of the historic range,
713 we will often not infer ancestors to be in the currently sparsely-occupied
714 portions. Other large scale population movements, such as one population
715 replacing another, may also partially obscure the geographic locations of an-
716 cestors. Over the past decade we have learned about numerous large-scale
717 movements of human populations alongside the expansions of archaeological
718 cultures, a fact fairly hidden from view by contemporary samples that only
719 ancient DNA could bring to light ([Slatkin and Racimo, 2016](#); [Reich, 2018](#)).
720 One obvious way to improve our method then, is to include ancient samples.
721 Given that it is now possible to include high-coverage ancient genomes in
722 tree sequences ([Speidel et al., 2021](#); [Wohns et al., 2021](#)), it is straightforward
723 to include this information in our likelihoods (ancient samples are treated as
724 any other, we calculate the shared times of these lineages with themselves
725 and with all other sample lineages), influencing both our dispersal estimates
726 and inferred ancestral locations. This should help, in particular, in locating
727 ancestors that are closely related to the ancient samples and for detecting
728 population-scale movements, such as range shifts, contractions, and replace-
729 ments.

730 Materials and Methods

731 Here we describe our methods to estimate dispersal rates and locate genetic
732 ancestors and how we applied these to simulations and *Arabidopsis thaliana*.

733 **The probability of the sampled locations given a tree
734 and the dispersal rate**

735 We first derive the probability of the sampled locations, \mathbf{L} , given a bifur-
736 cating tree topology (genealogy), \mathcal{G} , and associated branch lengths (times),
737 \mathcal{T} , which describe the coalescent history of the sample, and the dispersal
738 rate (covariance matrix), Σ . We derive this probability, $\mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}, \Sigma)$, com-
739 pounding the normally-distributed dispersal events each generation to give
740 Brownian motion down the tree in a similar manner to phylogenetic least
741 squares regression (Grafen, 1989; Harmon, 2019).

742 Let $\ell_i = [\ell_{i,1} \ \ell_{i,2} \ \dots \ \ell_{i,m}]^\top$ be the m -dimensional location of sample i ,
743 $\mathbf{L} = [\ell_1 \ \ell_2 \ \dots \ \ell_n]^\top$ the $n \times m$ matrix of spatial locations for all n samples,
744 and $\ell = [\ell_{1,1} \ \ell_{2,1} \ \dots \ \ell_{n,1} \ \ell_{1,2} \ \ell_{2,2} \ \dots \ \ell_{n,2} \ \dots \ \ell_{1,m} \ \ell_{2,m} \ \dots \ \ell_{n,m}]^\top$ the matrix of
745 locations represented as a vector of length nm . Let $\mathbf{S}_{\mathcal{G}, \mathcal{T}}$ be the $n \times n$ matrix
746 of shared evolutionary time (in generations) between each sample lineage in
747 the tree defined by \mathcal{G} and \mathcal{T} .

748 Then, assuming per generation dispersal is multivariate normal with mean
749 displacement $\mathbf{0}$ and covariance matrix Σ , the probability of the locations
750 given the tree is

$$\mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}, \Sigma) = \mathbb{P}(\ell|\mathbf{S}_{\mathcal{G}, \mathcal{T}}, \Sigma) \sim \mathcal{N}\left(\widehat{\mathbf{D}\ell_A}, \mathbf{S}_{\mathcal{G}, \mathcal{T}} \otimes \Sigma\right), \quad (1)$$

751 where

$$\widehat{\ell_A} = [(\mathbf{1}^\top \mathbf{S}_{\mathcal{G}, \mathcal{T}}^{-1} \mathbf{1})^{-1} (\mathbf{1}^\top \mathbf{S}_{\mathcal{G}, \mathcal{T}}^{-1} \mathbf{L})]^\top \quad (2)$$

752 is the maximum likelihood estimate for the location of the most recent com-
753 mon ancestor given the tree, $\mathbf{1}$ is a column vector of n ones, and \mathbf{D} is a
754 $nm \times m$ design matrix whose i, j^{th} entry is 1 if $(j-1)n < i \leq jn$ and 0
755 otherwise.

756 **Mean centering the locations**

757 We can remove any dependence on the (unknown) location of the most recent
758 common ancestor by mean centering the data. The mean centered locations
759 and associated matrix of shared times are $\mathbf{X} = \mathbf{ML}$ and $\mathbf{V}_{\mathcal{G}, \mathcal{T}} = \mathbf{MS}_{\mathcal{G}, \mathcal{T}}\mathbf{M}^\top$,
760 where \mathbf{M} is an $(n-1) \times n$ matrix with $n-1/n$ on the diagonal and $-1/n$
761 elsewhere. We only use those (sub)trees (see below) with $n > 1$ samples (i.e.,
762 trees with only one sample contain no information about the dispersal rate

763 if we do not know where the ancestor was). We then have

$$\mathbb{P}(\mathbf{X}|\mathcal{G}, \mathcal{T}, \Sigma) = \mathbb{P}(\mathbf{x}|\mathbf{V}_{\mathcal{G}, \mathcal{T}}, \Sigma) \sim \mathcal{N}(\mathbf{0}, \mathbf{V}_{\mathcal{G}, \mathcal{T}} \otimes \Sigma), \quad (3)$$

764 where \mathbf{x} is the vector representation of \mathbf{X} (as ℓ is to \mathbf{L}).

765 Chopping the tree into subtrees

766 We will usually want to only consider dispersal more recently than some cut-
767 off time, T , since deeper genealogical history may contain little geographic
768 information (Wilkins, 2004). When this is the case we cut the full tree of at T ,
769 leaving us with $n_T \in [1, n]$ subtrees, where all coalescent times are $< T$, and
770 we then calculate the shared times in each subtree independently. Letting
771 $\mathbb{P}(\mathbf{x}_i|\mathbf{V}_{\mathcal{G}_i, \mathcal{T}_i}, \Sigma)$ be the probability of the (mean centered and vectorized) loca-
772 tions of the samples in subtree i given the (mean centered) matrix of shared
773 evolutionary times in that subtree and the dispersal rate, the probability of
774 the locations of all samples is

$$\mathbb{P}(\mathbf{X}|\mathcal{G}, \mathcal{T}, \Sigma, T) = \prod_{i=1}^{n_T} \mathbb{P}(\mathbf{x}_i|\mathbf{V}_{\mathcal{G}_i, \mathcal{T}_i}, \Sigma) \sim \prod_{i=1}^{n_T} \mathcal{N}(\mathbf{0}, \mathbf{V}_{\mathcal{G}_i, \mathcal{T}_i} \otimes \Sigma). \quad (4)$$

775 This has the added benefit of expressing the probability as a function of
776 smaller submatrices of shared evolutionary times, which are faster to invert.

777 Dividing time into epochs

778 We can extend this model to allow dispersal to vary through time by dividing
779 time into epochs and assuming dispersal is only constant within each. If we
780 divide time into K epochs (as described by a vector of split times, \mathbf{t}), so that
781 the (mean centered) shared time matrix in the i^{th} epoch is $\mathbf{V}_{\mathcal{G}, \mathcal{T}}^{(i)}$, then

$$\begin{aligned} \mathbb{P}(\mathbf{X}|\mathcal{G}, \mathcal{T}, \Sigma_1, \dots, \Sigma_K, \mathbf{t}) &= \mathbb{P}(\mathbf{x}|\mathbf{V}_{\mathcal{G}, \mathcal{T}}^{(1)}, \dots, \mathbf{V}_{\mathcal{G}, \mathcal{T}}^{(E)}, \Sigma_1, \dots, \Sigma_E) \\ &\sim \mathcal{N}\left(\mathbf{0}, \sum_{i=1}^K (\mathbf{V}_{\mathcal{G}, \mathcal{T}}^{(i)} \otimes \Sigma_i)\right). \end{aligned} \quad (5)$$

782 Thus different dispersal epochs are easily accounted for.

783 **Importance sampling**

784 The calculations above, which give the likelihood of the dispersal rate given
 785 the sample locations, were all predicated on knowing the tree with certainty,
 786 which will not be the case when inferring trees from genetic data. Fur-
 787 ther, inferring a tree may involve assumptions (such as panmixia) that are
 788 inconsistent with the model we are using. To deal with the uncertainty
 789 and bias we calculate the likelihood of our parameters given the data using
 790 importance sampling, a likelihood ratio, and Monte Carlo approximation.
 791 Importance sampling corrects the expectation of the likelihood by reweight-
 792 ing draws from an ‘incorrect’ (proposal) distribution to match the ‘correct’
 793 (target) distribution. These importance weights downweight draws from the
 794 proposal distribution that are likely under the proposal distribution but un-
 795 likely under the target distribution and upweight draws that are likely in the
 796 target distribution and unlikely under the proposal distribution. Importance
 797 sampling techniques to sample genealogies in population genetics have been
 798 applied a number of settings as coalescent models provide convenient pri-
 799 ors on trees and it is often challenging to sample genealogies consistent with
 800 data (Griffiths and Tavaré, 1997; Stephens and Donnelly, 2000; Meligkotsidou
 801 and Fearnhead, 2007). See Stern et al. (2019, 2021) for recent applications of
 802 these ideas to marginal trees inferred along a recombining sequence, whose
 803 general approach we follow below.

804 The data we have are the locations of the samples, \mathbf{L} , and their hap-
 805 lotypes, \mathbf{H} . From this data we want to infer two unknowns, a tree topol-
 806 ogy, \mathcal{G} , and associated branch lengths, \mathcal{T} , and use these to estimate the
 807 likelihood the dispersal rate, Σ . Put another way, we want to estimate
 808 $\mathbb{E}_{\mathcal{G}, \mathcal{T} | \Sigma} [\mathbb{P}(\mathbf{L}, \mathbf{H} | \mathcal{G}, \mathcal{T}, \Sigma)]$, which is the likelihood of our parameters given the
 809 data, integrated over the unknowns.

810 Current methods to infer tree topologies and times along a recombining
 811 sequence (Rasmussen et al., 2014; Speidel et al., 2019; Kelleher et al., 2019;
 812 Wohns et al., 2021) assume panmictic, well-mixed populations. This implies
 813 we cannot sample topologies and branch lengths, \mathcal{G} and \mathcal{T} , under our spatial
 814 model, creating potential bias. To correct for this bias we use importance
 815 sampling,

$$\begin{aligned} L(\Sigma) &:= \mathbb{E}_{\mathcal{G}, \mathcal{T} | \Sigma} [\mathbb{P}(\mathbf{L}, \mathbf{H} | \mathcal{G}, \mathcal{T}, \Sigma)] \\ &= \mathbb{E}_{\mathcal{G}, \mathcal{T} | \mathbf{H}, \text{panmixia}} \left[\frac{\mathbb{P}(\mathbf{L}, \mathbf{H} | \mathcal{G}, \mathcal{T}, \Sigma) \mathbb{P}(\mathcal{G}, \mathcal{T} | \Sigma)}{\mathbb{P}(\mathcal{G}, \mathcal{T} | \mathbf{H}, \text{panmixia})} \right], \end{aligned} \quad (6)$$

816 where $\mathbb{P}(\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia})$ is the distribution of topologies and branch lengths
 817 we sample from.

818 The probabilities containing the genetic data, \mathbf{H} , are complicated to cal-
 819 culate. To simplify we divide by the likelihood of panmixia given the data,
 820 $L(\text{panmixia}) = \mathbb{P}(\mathbf{H}|\text{panmixia})$, to work with the likelihood ratio

$$\begin{aligned} \text{LR}(\boldsymbol{\Sigma}) &:= \frac{L(\boldsymbol{\Sigma})}{L(\text{panmixia})} \\ &= \mathbb{E}_{\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia}} \left[\frac{\mathbb{P}(\mathbf{L}, \mathbf{H}|\mathcal{G}, \mathcal{T}, \boldsymbol{\Sigma}) \mathbb{P}(\mathcal{G}, \mathcal{T}|\boldsymbol{\Sigma})}{\mathbb{P}(\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia}) \mathbb{P}(\mathbf{H}|\text{panmixia})} \right] \\ &= \mathbb{E}_{\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia}} \left[\frac{\mathbb{P}(\mathbf{L}, \mathbf{H}|\mathcal{G}, \mathcal{T}, \boldsymbol{\Sigma}) \mathbb{P}(\mathcal{G}, \mathcal{T}|\boldsymbol{\Sigma})}{\mathbb{P}(\mathbf{H}, \mathcal{G}, \mathcal{T}|\text{panmixia})} \right] \\ &= \mathbb{E}_{\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia}} \left[\frac{\mathbb{P}(\mathbf{H}|\mathcal{G}, \mathcal{T}) \mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}, \boldsymbol{\Sigma}) \mathbb{P}(\mathcal{G}, \mathcal{T}|\boldsymbol{\Sigma})}{\mathbb{P}(\mathbf{H}|\mathcal{G}, \mathcal{T}) \mathbb{P}(\mathcal{G}, \mathcal{T}|\text{panmixia})} \right] \\ &= \mathbb{E}_{\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia}} \left[\frac{\mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}, \boldsymbol{\Sigma}) \mathbb{P}(\mathcal{G}, \mathcal{T}|\boldsymbol{\Sigma})}{\mathbb{P}(\mathcal{G}, \mathcal{T}|\text{panmixia})} \right], \end{aligned} \tag{7}$$

821 where the third step assumes the genetic data (\mathbf{H}) is conditionally indepen-
 822 dent of the spatial parameters ($\boldsymbol{\Sigma}$ or panmixia) and locations (\mathbf{L}) given the
 823 tree (\mathcal{G} and \mathcal{T}). We can approximate this expectation using Monte Carlo
 824 sampling

$$\widehat{\text{LR}(\boldsymbol{\Sigma})} = \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\mathbf{L}|\mathcal{G}_m, \mathcal{T}_m, \boldsymbol{\Sigma}) \mathbb{P}(\mathcal{G}_m, \mathcal{T}_m|\boldsymbol{\Sigma})}{\mathbb{P}(\mathcal{G}_m, \mathcal{T}_m|\text{panmixia})}, \tag{8}$$

825 where $\mathcal{G}_m, \mathcal{T}_m \sim \mathbb{P}(\mathcal{G}, \mathcal{T}|\mathbf{H}, \text{panmixia})$ is sampled using Markov chain Monte
 826 Carlo.

827 We make two final simplifications. First, we will use a model of branch-
 828 ing Brownian motion for $\mathbb{P}(\mathcal{G}, \mathcal{T}|\boldsymbol{\Sigma})$ and the standard neutral coalescent for
 829 $\mathbb{P}(\mathcal{G}, \mathcal{T}|\text{panmixia})$. Under both of these models the probability of the topol-
 830 ogy, $\mathbb{P}(\mathcal{G}|\text{panmixia})$, is equivalent (and uniform). Second, we will use **Relate**
 831 ([Speidel et al., 2019](#)) to infer topologies and branch lengths. **Relate** returns
 832 a single topology and allows resampling over branch lengths conditional on
 833 this topology. We therefore take the topology as given and integrate only

834 over branch lengths. Putting these two simplifications together,

$$\begin{aligned}
\widehat{\text{LR}(\Sigma)} &= \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\mathbf{L}|\mathcal{G}_m, \mathcal{T}_m, \Sigma) \mathbb{P}(\mathcal{G}_m, \mathcal{T}_m|\Sigma)}{\mathbb{P}(\mathcal{G}_m, \mathcal{T}_m|\text{panmixia})} \\
&= \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\mathbf{L}|\mathcal{G}_m, \mathcal{T}_m, \Sigma) \mathbb{P}(\mathcal{T}_m|\mathcal{G}_m, \Sigma) \mathbb{P}(\mathcal{G}_m|\Sigma)}{\mathbb{P}(\mathcal{T}_m|\mathcal{G}_m, \text{panmixia}) \mathbb{P}(\mathcal{G}_m|\text{panmixia})} \\
&= \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\mathbf{L}|\mathcal{G}_m, \mathcal{T}_m, \Sigma) \mathbb{P}(\mathcal{T}_m|\mathcal{G}_m, \Sigma)}{\mathbb{P}(\mathcal{T}_m|\mathcal{G}_m, \text{panmixia})} \\
&\approx \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}_m, \Sigma) \mathbb{P}(\mathcal{T}_m|\mathcal{G}, \Sigma)}{\mathbb{P}(\mathcal{T}_m|\mathcal{G}, \text{panmixia})},
\end{aligned} \tag{9}$$

835 where $\mathcal{T}_m \sim \mathbb{P}(\mathcal{T}|\mathbf{H}, \text{panmixia}, \mathcal{G})$ is sampled using Markov chain Monte
836 Carlo. Note that because the probability of a topology is the same in the two
837 models, and therefore cancels out, our method can immediately be extended
838 to integrate over both topologies and branch lengths.

839 To find the maximum likelihood estimate of Σ we numerically search for
840 the Σ that maximizes this likelihood ratio. We measure the variance around
841 this estimate with the Hessian matrix of the likelihood surface returned by
842 the numerical search. To use information from multiple loci, e.g., to estimate
843 a per-chromosome or genome-wide dispersal parameter, we can multiply the
844 likelihood ratios together to give a composite likelihood ratio. This is a
845 composite likelihood because it ignores correlations in the trees between loci
846 (Hudson, 2001; Larribe and Fearnhead, 2011; Varin et al., 2011). The max-
847 imum likelihood of genome-wide parameters from composite likelihood are
848 known to be statistically consistent in the presence of such correlations (Wiuf,
849 2006), but the likelihood surface is overly peaked due to ignoring the corre-
850 lations in the data. However, in practice we take loci that are far enough
851 apart that their trees are only weakly correlated more than a few tens of
852 generations back.

853 We next show how we derive the probability distributions of the branch
854 lengths in the tree that we use as the weights in the importance sampler,
855 $\mathbb{P}(\mathcal{T}|\mathcal{G}, \text{panmixia})$ and $\mathbb{P}(\mathcal{T}|\mathcal{G}, \Sigma)$.

856 **The probability of the coalescence times given the tree topology
857 and panmixia**

858 Relate and other tree construction methods infer the times in the trees un-
859 der a model of panmixia with variable population size through time (with
860 this demographic history inferred genomes as part of the method). The
861 probability of the coalescence times given the tree topology and panmixia,
862 $\mathbb{P}(\mathcal{T}|\mathcal{G}, \text{panmixia})$, can be derived from the standard neutral coalescent al-
863 lowing for (stepwise) changing effective population size ([Griffiths and Tavaré, 1994; Meligkotsidou and Fearnhead, 2007](#)).

865 Let the number of samples be n and let \mathbf{U} be the $n - 1$ coalescence times,
866 ordered from most recent to most historic. Then

$$\mathbb{P}(\mathcal{T}|\mathcal{G}, \text{panmixia}) = \prod_{i=1}^{n-1} \mathbb{P}(U_i = u_i | U_{i-1} = u_{i-1}), \quad (10)$$

867 where

$$\begin{aligned} & \mathbb{P}(U_i = u_i | U_{i-1} = u_{i-1}) \\ &= \binom{n - (i-1)}{2} \frac{1}{2N(u_i)} \exp \left[-\binom{n - (i-1)}{2} [\Lambda(u_i) - \Lambda(u_{i-1})] \right] \end{aligned} \quad (11)$$

868 with

$$\Lambda(u) = \int_0^u \frac{1}{2N(t)} dt, \quad (12)$$

869 where $N(t)$ is the effective population size at time t in the past.

870 This says that the coalescence times are produced by a Markov process, so
871 that the probability of the i^{th} coalescence being at time u_i is independent of \mathbf{U}
872 once conditioned on the time of the previous coalescence, u_{i-1} . Between time
873 u_{i-1} and u_i we have $n - (i-1)$ lineages. The probability of coalescence at time
874 u_i is then approximately distributed as a time-inhomogeneous exponential
875 random variable with instantaneous rate $\lambda(u) = \binom{n - (i-1)}{2} / [2N(u)]$.

876 When effective population size is piecewise constant, split into K epochs
877 with effective population size N_k between time τ_k and τ_{k+1} , then (see Ap-
878 pendix A of [Stern et al., 2021](#))

$$\Lambda(u) = \left[\sum_{k=0}^{b(u)} \frac{1}{2N_k} (\tau_{k+1} - \tau_k) \right] + \frac{1}{2N_{b(u)+1}} (u - \tau_{b(u)+1}), \quad (13)$$

879 where $b(u) = \max(k \in \{0, 1, 2, \dots, K - 1\} : u > \tau_{k+1})$ is the last epoch to end
 880 before time u .

881 If we only want the probability of coalescence times back as far as time
 882 T , we consider only these $m \leq n - 1$ coalescence times and multiply the
 883 probability of these times by the probability of no coalescence between time
 884 u_m and T ,

$$\begin{aligned} & \mathbb{P}(\mathcal{T}|\mathcal{G}, \text{panmixia}, T) \\ &= \exp \left[- \binom{n-m}{2} [\Lambda(T) - \Lambda(u_m)] \right] \prod_{i=1}^m \mathbb{P}(U_i = u_i | U_{i-1} = u_{i-1}). \end{aligned} \quad (14)$$

885 **The probability of the coalescence times given the tree topology
 886 and spatial model**

887 A number of different fully spatial models of the coalescent have been de-
 888 veloped. However, for our purposes they are computational challenging to
 889 work with and so we pursue a simpler analytically-tractable approximation
 890 called branching Brownian motion (BBM), or the Brownian-Yule process.
 891 Such approximations have been used previously to approximate the short
 892 timescale dynamics of genealogies ([Edwards, 1970](#); [Rannala and Yang, 1996](#);
 893 [Meligkotsidou and Fearnhead, 2007](#)). We view the rate of branching (λ) in
 894 the BBM process as a nuisance parameter, but informally it is proportional
 895 to the inverse of the local population density and so it might be a useful
 896 parameter to explore for future work (for a recent application see [Ringbauer
 et al., 2017](#)).

897 Let T be the time in the past we wish to start the process (which can be
 898 no greater than the time to the most recent common ancestor in the tree).
 899 Let n_0 be the number of ancestral lineages at this time (which can be no
 900 smaller than 2). Taking a forward in time perspective, let \mathbf{U}' be the birth
 901 times, $U'_i = T - U_{n-i}$, i.e., U'_i is the i^{th} split time, measured forward in
 902 time from T , transitioning from i to $i + 1$ lineages. Then under a pure birth
 903 (Yule) process with per capita birth rate λ the joint probability density of
 904 the $n - n_0$ split times and ending up with n samples (i.e., no birth once n

906 lineages existed) is

$$\begin{aligned} f(\mathbf{u}', n | \lambda, n_0, T) &= \left(\prod_{i=n_0}^{n-1} \lambda i \exp[-\lambda i(u'_i - u'_{i-1})] \right) \exp[-\lambda n(T - u'_{n-1})] \\ &= \lambda^{n-n_0} \frac{(n-1)!}{(n_0-1)!} \exp \left[-\lambda \left(nT + \sum_{i=n_0}^{n-1} u'_i \right) \right], \end{aligned} \quad (15)$$

907 where we have set $u'_{n_0-1} = 0$. Now, the probability that n_0 lineages pro-
908 duce exactly n lineages in time T under a pure birth process is given by a
909 negative binomial with $n - n_0$ successes, n_0 failures, and success probability
910 $1 - \exp(-\lambda T)$,

$$\mathbb{P}(n | \lambda, n_0, T) = \binom{n-1}{n-n_0} \exp(-\lambda T n_0) [1 - \exp(-\lambda T)]^{n-n_0}. \quad (16)$$

911 The probability of the coalescence times under the Yule process, which we
912 take to be the probability of the times given the topology and spatial model,
913 is therefore proportional to the ratio of Equations (15) and (16),

$$\begin{aligned} \mathbb{P}(T | G, \Sigma) &= \mathbb{P}(\mathbf{U}' = \mathbf{u}' | \lambda, n_0, T, n) \\ &\propto \frac{f(\mathbf{u}', n | \lambda, n_0, T)}{\mathbb{P}(n | \lambda, n_0, T)} \\ &= (n - n_0)! \left(\frac{\lambda \exp(-\lambda T)}{1 - \exp(-\lambda T)} \right)^{n-n_0} \exp \left(-\lambda \sum_{i=n_0}^{n-1} u'_i \right). \end{aligned} \quad (17)$$

914 See Edwards (1970); Rannala and Yang (1996); Meligkotsidou and Fearnhead
915 (2007) for more on the Yule process.

916 Conditioning on sampling locations

917 As many studies choose sampling locations before choosing samples, we would
918 like to condition on sampling locations in our inferences of dispersal rate,
919 i.e., use the conditional likelihood $\mathbb{P}(\mathbf{L} | \mathcal{G}, \mathcal{T}, \Sigma) / \mathbb{P}(\mathbf{L} | \Sigma)$. The denomina-
920 tor, $\mathbb{P}(\mathbf{L} | \Sigma)$, is the probability of the locations of the descendent tips of a
921 branching Brownian motion given the dispersal rate. This is the numerator
922 $\mathbb{P}(\mathbf{L} | \mathcal{G}, \mathcal{T}, \Sigma)$ after the genealogy and branch times have been integrated out.
923 Calculating the integrals over genealogies and branch times could be achieved

by Monte Carlo simulation, but would be computationally intensive. However, as shown by (Meligkotsidou and Fearnhead, 2007), in one dimension (say x) the probability of the locations, $\mathbb{P}(\mathbf{L}|\sigma_x)$, is independent of dispersal, σ_x , if a scale invariant prior is placed on the branching rate. Intuitively, this follows from the fact that if we do not *a priori* know the scale of the coalescent rate then we do not know *a priori* how long the branch lengths are, and so we do not know how fast lineages need to disperse to get to where they are today. Therefore, in that case, conditioning on the locations does not change the likelihood surface of σ_x and so can be ignored in the inference of dispersal. This result also holds in two dimensions if we constrain $\sigma_x = \sigma_y$. Sadly it does not hold in general for arbitrary dispersal matrix Σ . For example, if our samples (the end points of BBM) were distributed widely along the x axis but varied little in displacement along the y axis, this would be more likely if $\sigma_y/\sigma_x \ll 1$. Thus, while the overall magnitude of dispersal is not affected by conditioning on the sampling locations, the anisotropy of dispersal may be.

As one of our results is that rates of longitudinal dispersal are higher than latitudinal rates, we wanted to make sure that the lack of conditioning on sampling location did not drive this result. To do this we ran our method looking at only longitudinal dispersal, ignoring the samples latitudes, and then again looking at only latitudinal dispersal, ignoring the samples longitudes. This reduces the problem back to one dimension where the conditioning does not matter (as discussed above and in Meligkotsidou and Fearnhead, 2007). Doing this we find essentially identical dispersal rates as before, justifying our result of a larger longitudinal dispersal rate despite not conditioning on the sample locations.

Maximum likelihood estimates of the dispersal rate given a tree

If we take the tree as fixed, Equation (1) gives the likelihood of the dispersal rate given the locations. This implies the maximum likelihood estimate (MLE) of the dispersal rate given the tree is

$$\widehat{\Sigma} = \frac{(\boldsymbol{\ell} - \widehat{\boldsymbol{\ell}}_A \mathbf{1}^\top)^\top \mathbf{S}_{G,T}^{-1} (\boldsymbol{\ell} - \widehat{\boldsymbol{\ell}}_A \mathbf{1}^\top)}{n}. \quad (18)$$

955 This estimate is, however, biased. Using Equation (3), the restricted maxi-
 956 mum likelihood estimate of the dispersal rate given the tree is

$$\widehat{\Sigma}^* = \frac{\mathbf{x}^\top \mathbf{V}_{\mathcal{G},\mathcal{T}}^{-1} \mathbf{x}}{n-1}, \quad (19)$$

957 which is unbiased and equivalent to $n\widehat{\Sigma}/(n-1)$.

958 In practice, when the tree has been chopped into subtrees, dispersal is
 959 allowed to vary through time, or we importance sample then we do not have
 960 an analytical expression for the MLE dispersal rate and instead numerically
 961 search for the dispersal rate which maximizes the likelihood (Equations (4),
 962 (5), (9)).

963 Likelihood of the location of a genetic ancestor given 964 the tree

965 A point on a branch of a tree represents a genetic ancestor. Choosing one
 966 such point, if the ancestor's lineage shares \mathbf{s}_a evolutionary time the sample
 967 lineages and s_a time with itself, then, by the properties of the conditional
 968 normal distribution, the probability that the ancestor was at location ℓ_a is

$$\mathbb{P}(\ell_a | \ell, \mathbf{S}_{\mathcal{G},\mathcal{T}}, \Sigma, \mathbf{s}_a) \sim \mathcal{N}\left(\widehat{\ell}_a, \widehat{\mathbf{S}}\right), \quad (20)$$

969 where

$$\widehat{\ell}_a = \widehat{\ell}_A + \mathbf{s}_a^\top \mathbf{S}_{\mathcal{G},\mathcal{T}}^{-1} (\ell - \widehat{\ell}_A \mathbf{1}^\top) \quad (21)$$

970 is the maximum likelihood estimate of the ancestor's location and

$$\widehat{\mathbf{S}} = (s_a - \mathbf{s}_a^\top \mathbf{S}_{\mathcal{G},\mathcal{T}}^{-1} \mathbf{s}_a) \Sigma \quad (22)$$

971 is the uncertainty (covariance) in the maximum likelihood estimate. Here
 972 $\mathbf{S}_{\mathcal{G},\mathcal{T}}^{-1}$ is the generalized inverse of $\mathbf{S}_{\mathcal{G},\mathcal{T}}$.

973 This approach gives the correct maximum likelihood estimate (Equation
 974 (21)) but the uncertainty (Equation (22)) is incorrect because the approach
 975 implicitly assumes we know with certainty where the most recent common
 976 ancestor was.

977 The correct uncertainty is derived by mean centering. The mean centered
 978 matrices of sample locations, \mathbf{X} , and shared times among the samples, $\mathbf{V}_{\mathcal{G},\mathcal{T}}$
 979 are derived as above. The mean centered vector of shared times between

980 the ancestor and the samples is $\mathbf{v}_a = \mathbf{M}(\mathbf{s}_a - \mathbf{S}_{\mathcal{G},\mathcal{T}}\mathbf{1}/n)$, with M and n
 981 the mean centering matrix and number of samples, as above. The mean
 982 centered shared time of the ancestor with itself is $v_a = \mathbf{m}^\top \mathbf{S}_{\mathcal{G},\mathcal{T},a} \mathbf{m}$, where \mathbf{m}
 983 is a $(n+1)$ -column vector with 1 as the first entry and $-1/n$ elsewhere and
 984 $\mathbf{S}_{\mathcal{G},\mathcal{T},a} = \begin{pmatrix} s_a & \mathbf{s}_a^\top \\ \mathbf{s}_a & \mathbf{S}_{\mathcal{G},\mathcal{T}} \end{pmatrix}$.

985 The probability the ancestor was at location ℓ_a is then

$$\mathbb{P}(\ell_a | \mathbf{x}, \mathbf{V}_{\mathcal{G},\mathcal{T}}, \Sigma, \mathbf{v}_a, v_a) \sim \mathcal{N}\left(\hat{\ell}_a, \hat{\mathbf{S}}^*\right), \quad (23)$$

986 where the maximum likelihood estimate has not changed but can be rewritten

$$\hat{\ell}_a = \bar{\mathbf{L}} + \mathbf{v}_a^\top \mathbf{V}_{\mathcal{G},\mathcal{T}}^{-1} \mathbf{x}, \quad (24)$$

987 with $\bar{\mathbf{L}}$ the mean location of the samples. The uncertainty in this estimate is

$$\hat{\mathbf{S}}^* = (v_a - \mathbf{v}_a^\top \mathbf{V}_{\mathcal{G},\mathcal{T}}^{-1} \mathbf{v}_a) \Sigma. \quad (25)$$

988 This uncertainty is the correct uncertainty. For example, it produces a linear
 989 increase in the variance of the estimate as we move up from a sample when
 990 there is only a single sample, i.e., this is simple Brownian motion. When there
 991 are two samples with a recent common ancestor at time T , the variance in
 992 ancestor location along this tree back to the most recent common ancestor is
 993 maximized at the most recent common ancestor, $T/2$, i.e., this is a Brownian
 994 bridge. And so on.

995 As with the dispersal likelihood, to reduce dependencies on distant times
 996 we can chop the tree at some time T and use only the subtree containing the
 997 ancestor of interest. Often the most recent common ancestor of the resulting
 998 subtree occurs more recently than T but we want to infer the location all the
 999 way back to T , and so we need to infer the location of the common ances-
 1000 tral lineage past the most recent common ancestor. However, this method
 1001 extends naturally to times beyond the most recent common ancestor, implic-
 1002 itly modelling simple Brownian motion (a fixed mean and linearly increasing
 1003 variance) up the stem of the tree. Similarly, the ancestor need not have
 1004 descendants in the sample – our method implicitly adds simple Brownian
 1005 motion down the branch leading to the ‘hanging’ ancestor (which is useful
 1006 when we want to ignore the location of a sample and locate it or its ancestors
 1007 using the remaining sample locations and the trees). When there is only a
 1008 single sample we cannot mean center and instead model simple Brownian
 1009 motion up from the sample.

1010 We want to allow dispersal rates to vary across epochs. To account for
 1011 this in locating ancestors, instead of simply calculating the shared evolutionary
 1012 times between lineages (and then multiplying by a constant dispersal
 1013 rate) we must explicitly calculate the covariance between each lineage. The
 1014 covariance between two lineages in any one epoch is the (Kronecker) product
 1015 of their shared time in that epoch and the dispersal rate in that epoch. The
 1016 total covariance between the lineages is then the sum of these covariances
 1017 over epochs. We can then follow the approach above to get the probability
 1018 distribution of ancestor locations (Equation (23)) after replacing \mathbf{M} with
 1019 $\mathbf{M} \otimes \mathbf{I}_d$, \mathbf{m} with $\mathbf{m} \otimes \mathbf{I}_d$, $\mathbf{1}$ with $\mathbf{1} \otimes \mathbf{I}_d$, \mathbf{x} (in Equation (24)) with $\mathbf{x} \otimes \mathbf{1}_d$,
 1020 and Σ (in Equation (25)) with 1, where \mathbf{I}_d is the identity matrix and $\mathbf{1}_d$ is a
 1021 column vector of 1s, each with the same dimension (d) as Σ .

1022 Finally, to integrate over uncertainty and reduce bias we also want to use
 1023 importance sampling. We do this by replacing $\mathbb{P}(\mathbf{L}|\mathcal{G}, \mathcal{T}, \Sigma)$ in Equation (9)
 1024 with Equation (23),

$$\widehat{\text{LR}}(\boldsymbol{\ell}_a) = \frac{1}{M} \sum_{m=1}^M \frac{\mathbb{P}(\boldsymbol{\ell}_a | \mathbf{x}, \mathbf{V}_{\mathcal{G}, \mathcal{T}_m}, \Sigma, \mathbf{v}_a, v_a) \mathbb{P}(\mathcal{T}_m | \mathcal{G}, \Sigma)}{\mathbb{P}(\mathcal{T}_m | \mathcal{G}, \text{panmixia})}, \quad (26)$$

1025 and generally use the maximum composite likelihood estimates of Σ and λ
 1026 (e.g., per-chromosome or genome-wide estimates). We measure the variance
 1027 around this estimate with the Hessian matrix returned by the numerical
 1028 search.

1029 Best Linear Unbiased Predictions of Locations.

1030 The above approach requires a numerical search for the maximum likelihood
 1031 ancestor location because a sum of normal distributions (Equation (26)) does
 1032 not, in general, follow a tractable distribution. An alternative approach,
 1033 however, is to calculate the maximum likelihood ancestor location for each
 1034 sampled tree (Equation (24)) and importance sample over these estimates.
 1035 Writing the maximum likelihood ancestor location given a sampled tree as
 1036 $\widehat{\boldsymbol{\ell}}_a(\mathcal{T}_m)$, we then have another estimate of the ancestor's location,

$$\widehat{\boldsymbol{\ell}}_a^* = \frac{1}{M} \sum_{m=1}^M \frac{\widehat{\boldsymbol{\ell}}_a(\mathcal{T}_m) \mathbb{P}(\mathcal{T}_m | \mathcal{G}, \Sigma)}{\mathbb{P}(\mathcal{T}_m | \mathcal{G}, \text{panmixia})}. \quad (27)$$

1037 Equation (27) is a ‘Best Linear Unbiased Predictor’ (BLUP) averaged over
 1038 the importance weights. We can calculate this directly and it is therefore, in

principle, faster to compute than the search over Equation (26). In practice we find that the MLE and weighted-BLUP estimators are very correlated. Throughout the paper we use the MLEs but both are implemented in our software.

Simulations

Spatially-explicit forward-time simulations

We performed simulations in SLiM v3.6 (Haller and Messer, 2019) with tree sequence recording (Haller et al., 2019). The SLiM code was adapted from the `pyslim` spatial vignette (https://pyslim.readthedocs.io/en/latest/vignette_space.html) to model non-overlapping generations.

Individuals are diploid for a single chromosome with $L = 10^8$ basepairs and per basepair recombination rate $r = 5 \times 10^{-9}$. Individuals exist in a two dimensional (2D) habitat, a square of width $W = 50$ with reflecting boundaries. Each generation begins with reproduction. Each individual acts once as a ‘mother’ and chooses a ‘father’ at random (individuals are hermaphrodites), weighted by their mating weights. The mating weight of an individual $d < 3\sigma$ distance from a mother follows a 2D normal distribution centered on the mother with variance σ^2 in both directions and no covariance. Individuals further than 3σ distance apart are ignored for efficiency (creating an Allee effect – i.e., a mother is not guaranteed to find a mate). Local density-dependence is modelled through competitive effects on fecundity. The strength of competitive interaction between two individuals $d < 3\sigma_c$ distance apart follows a 2D normal distribution centered on one of the individuals with variance $\sigma_c^2 = 0.5^2$ in both directions and no covariance (we again ignore more distant individuals). The number of offspring produced by a mating is Poisson, with mean $R/(1 + C/K)$, where C is the sum of interaction strengths the mother experiences, $R = 2$ is the mean number of offspring in the absence of competition, and $K = 2$ is the local carrying capacity. Each offspring disperses from its mother by a random 2D normal deviate with variance σ^2 in each dimension and no covariance. After reproduction all parents die and all offspring become adults in the next generation. We begin the population with $N_0 = W^2K = 5 \times 10^3$ individuals distributed uniformly at random across space. We end the simulation, and output the tree sequence, after $4N_0 = 2 \times 10^4$ generations.

1073 **True tree sequence of the sample**

1074 Using `pyslim` v0.6, we load the tree sequence into Python, sample $k/2 = 50$
1075 present-day individuals ($k = 100$ chromosomes) at random, and simplify
1076 the tree sequence to lineages ancestral to the sample. We next use `pyslim`
1077 to ‘recapitiate’ under the standard coalescent with recombination (Hudson,
1078 2002), with effective population size N_0 , to ensure all sampled lineages have
1079 coalesced. We then use `msprime` v1.0.1 (Kelleher et al., 2016b) to layer
1080 on neutral mutations with per basepair per generation mutation rate $U =$
1081 1.25×10^{-8} . This tree sequence represents the true genealogical history of
1082 the sample.

1083 **Inferred tree sequence of the sample**

1084 Because we will not know the true tree sequence of any natural sample, we
1085 also write this true tree sequence out as a VCF and use `Relate` v1.1.4 (Speidel
1086 et al., 2019) to infer the tree sequence. We give `Relate` the true uniform
1087 recombination map and mutation rate and an effective population size of
1088 $\theta/(4U)$, where θ is the observed mean genetic diversity (calculated from the
1089 tree sequence with `tskit` v0.3.5; Kelleher et al., 2018; Ralph et al., 2020). We
1090 then feed the resulting output to `Relate`’s ‘EstimatePopulationSize’ function
1091 to iteratively estimate a piece-wise constant effective population size and
1092 branch lengths, using 5 iterations and keeping all of the trees. We then
1093 use `Relate` to convert the anc/mut format to a `tskit` tree sequence for
1094 downstream processing.

1095 **True locations of ancestors**

1096 To analyze our ability to locate ancestors we also ran some `SLiM` simulations
1097 with `initializeTreeSeq(retainCoalescentOnly=F)` and
1098 `treeSeqRememberIndividuals(permanent=F)` options, meaning we remem-
1099 ber all individuals (including their location data) that are ancestors of the fi-
1100 nal population. When simplifying the tree sequence we use the `keep_unary=True`
1101 option, to keep all nodes (and their locations) that are ancestral to the sample
1102 (rather than just the coalescent nodes).

1103 **Processing the tree sequence**

1104 Each simplified tree sequence contains on the order of 10^4 trees. Here we
1105 use only 10^2 of these trees, uniformly sampled from the tree sequence, so
1106 that each is relatively independent of the others. For each of these trees we
1107 use `Relate`'s 'SampleBranchLengths' function to sample branch lengths from
1108 the posterior $M = 10$ times. We load resulting newick trees with `DendroPy`
1109 v4.5.2 ([Sukumaran and Holder, 2010](#)) and process each sampled tree. In
1110 processing, we first chop the trees off at T generations to create subtrees
1111 and for each subtree record the coalescence times, the probability of these
1112 coalescence times under the neutral coalescent (given the estimated piecewise
1113 constant population size from `Relate`), and the shared evolutionary times
1114 between each sample. This pre-processing speeds up the numerical search
1115 for maximum likelihood parameters.

1116 **Dispersal estimates**

1117 We approximate the maximum likelihood estimate (MLE) of dispersal rate
1118 at each tree by numerically searching for the maximum of Equation (9).
1119 We use the L-BFGS-B method of `SciPy` v1.6.2 ([Virtanen et al., 2020](#)) to
1120 numerically find the MLEs, which allows us to put bounds on the parameters.
1121 We search for the MLEs in terms of the standard deviation of dispersal in
1122 x (latitude) and y (longitude) (with a lower bound of 10^{-6} to prevent non-
1123 positive estimates) and the correlation between these two axes (with a lower
1124 bound of -0.99 and upper bound of 0.99 to prevent estimates with absolute
1125 value greater than 1). In the process we also find the MLE of the birth
1126 rate in the branching process (with a lower bound of 10^{-6} to prevent non-
1127 positive estimates). As the likelihood was much more sensitive to a given
1128 change in birth rate than the other parameters, we instead search for 10^2
1129 times the MLE birth rate and then rescale back to the original parameters,
1130 which makes the search for the MLE more efficient. We find the maximum
1131 composite likelihood estimate of dispersal and branching rate by searching
1132 for the parameters that maximize the product of the likelihoods over multiple
1133 loci.

1134 **Locating genetic ancestors**

1135 To locate a genetic ancestor at a particular locus and time we numerically
1136 search for the location that maximizes Equation (26), following the same

1137 approach as above with bounds of $(0, W)$ in both dimensions.

1138 *Arabidopsis thaliana* data

1139 Inferring the tree sequence

1140 We first downloaded the VCF of SNPs for 1135 *Arabidopsis thaliana* individuals from <https://1001genomes.org/data/GMI-MPI/releases/v3.1/> (Alonso-
1141 Blanco et al., 2016)) and used PLINK 2.0 (www.cog-genomics.org/plink/
1142 2.0/; Chang et al., 2015) to find SNPs with minor allele frequency < 0.05 .
1143 We then converted the matrix of imputed SNPs (https://1001genomes.org/data/GMI-MPI/releases/v3.1/SNP_matrix_imputed_hdf5) into a hap-
1144 loid ‘haps’ file (as required by `Relate`) for each chromosome, while filtering
1145 out those SNPs with minor allele frequency < 0.05 . A multi-species alignment
1146 for *Arabidopsis thaliana*, *Boechera stricta*, *Arabidopsis lyrata*, and *Mal-*
1147 *comia maritima* was kindly provided by Tyler Kent (University of Toronto).
1148 We used `bx-python` v0.8.9 to create a FASTA file containing the three out-
1149 groups from this alignment and, combined with the haps file described above,
1150 created the input file required by `est-sfs` (Keightley and Jackson, 2018),
1151 again for each chromosome separately. We then ran `est-sfs`, which in-
1152 corporates both the within-population polymorphism data and outgroup
1153 sequences, on each chromosome separately using the Kimura 2-parameter
1154 model (Kimura, 1980) to get the probability that each reference allele is
1155 ancestral. We then used `bx-python` to extract the reference sequence for
1156 *Arabidopsis thaliana* from the alignment and created the ancestral chromo-
1157 somes by making the alternate allele ancestral whenever the probability the
1158 reference allele was ancestral was < 0.5 . We then used this ancestral genome
1159 to create polarized haps files for each chromosome with `Relate`’s ‘FlipHap-
1160 sUsingAncestor’ function. The recombination map for each chromosome
1161 was downloaded from <https://www.eeb.ucla.edu/Faculty/Lohmueller/data/uploads/> (Salomé et al., 2012). The polarized haps files and recombi-
1162 nation maps were then used to infer the tree sequence with `Relate`, using per
1163 base pair per generation mutation rate $U = 7 \times 10^{-9}$ (Adrion et al., 2020) and
1164 haploid effective population size $2N_e = 1.7 \times 10^5$ (estimated from nucleotide
1165 diversity, $\pi = 4N_eU$, calculated with `tskit` `diversity()` statistic; Ralph
1166 et al., 2020). We then fed this output to `Relate`’s ‘EstimatePopulation-
1167 Size’ function (with some customizing of the script to allow it to work with
1168 haploids and spit out the anc/mut files) to iteratively estimate a piece-wise
1169
1170
1171

1172 constant effective population size and branch lengths, assuming $U = 7 \times 10^{-9}$
1173 and a single panmictic population, using 5 iterations and dropping half of
1174 the trees (the 50% with fewest mutations). We then converted the output to
1175 a `tskit` tree sequence using `Relate`'s 'ConvertToTreeSequence' function.

1176 **Nearly-identical samples**

1177 The 1001 Genomes dataset contains a number of nearly-identical samples as
1178 the result of selfing. A number of nearly identical pairs are separated by
1179 $> 1\text{km}$ and thus may represent long-distance migration or mis-assignment
1180 or mix-ups (Alonso-Blanco et al., 2016). We reduce the effect of these
1181 near identical samples by first calculating all pairwise genetic distances us-
1182 ing PLINK 1.9 (with the `--distance allele-ct` option; see Figure 3A in
1183 Alonso-Blanco et al., 2016). We then ignore all samples that differ at less
1184 than 10^3 base pairs from any other sample.

1185 **Removing dispersal outliers and samples without locations**

1186 We downloaded the metadata associated with the 1001 Genomes samples
1187 from https://raw.githubusercontent.com/hagax8/arabidopsis_viz/master/
1188 `data/dataframe_1001G.csv` (see also <https://1001genomes.org/acccessions.html>), which includes location data and previous admixture assignments
1189 (Alonso-Blanco et al., 2016) (see <https://1001genomes.github.io/admixture-map/>)
1190 for ancestry proportions). We remove the outlier samples in North America
1191 ($n = 125$) and Japan ($n = 2$), as well as those without locations ($n = 4$),
1192 from the dispersal rate and genetic ancestor location likelihoods (they remain
1193 in the importance sample weights).

1195 **Dispersal estimates**

1196 We estimate the MLE of dispersal rate at each chosen tree as described
1197 above for simulations. To reduce computation time we search for maximum
1198 composite likelihoods across trees for each chromosome separately, rather
1199 than genome-wide.

1200 We convert the estimates of dispersal rate (standard deviations) from de-
1201 grees to kilometres by multiplying the latitudinal estimates by $\cos(x\pi/180)/111$,
1202 where x is the mean latitude of all (filtered) samples, and multiplying the
1203 longitudinal estimates by 110.

1204 **Locating genetic ancestors**

1205 We locate ancestors as described above for simulations, using bounds of
1206 $(-90, 90)$ for latitude and $(-180, 180)$ for longitude.

1207 **Data availability statement**

1208 All code used to perform the analyses in this study can be found at [\[link\]](#). A
1209 Python package of our method is available at [\[link\]](#)

1210 **Acknowledgements**

1211 We thank Tyler Kent for the multispecies genome alignments, Yan Wong,
1212 Ben Haller, and Peter Ralph for the `retainCoalescentOnly` and `permanent`
1213 updates to `SLiM/pyslim`, Aaron Stern for making his importance sampling
1214 code freely available, Leo Speidel for helpful discussion regarding `Relate`,
1215 Vince Buffalo for the introduction to `Snakemake`, and Nick Barton for the
1216 healthy skepticism about per-locus estimates of ancestor locations at deep
1217 times. Funding provided by Banting (Canada) and Center for Population
1218 Biology (UC Davis) fellowships (awarded to MMO) and the National Institute
1219 of General Medical Sciences of the National Institutes of Health (NIH
1220 R01 GM108779 and R35 GM136290, awarded to GC). Computations were
1221 performed on the Niagara supercomputer at the SciNet HPC Consortium.
1222 SciNet is funded by: the Canada Foundation for Innovation; the Government
1223 of Ontario; Ontario Research Fund - Research Excellence; and the University
1224 of Toronto.

1225 **Competing interests**

1226 We have no financial and non-financial competing interests to declare.

1227 **References**

1228 Adriion, J. R., Cole, C. B., Dukler, N., Galloway, J. G., Gladstein, A. L.,
1229 Gower, G., Kyriazis, C. C., Ragsdale, A. P., Tsambos, G., Baumdicker,
1230 F., Carlson, J., Cartwright, R. A., Durvasula, A., Gronau, I., Kim,

- 1231 B. Y., McKenzie, P., Messer, P. W., Noskova, E., Ortega-Del Vecchyo,
1232 D., Racimo, F., Struck, T. J., Gravel, S., Gutenkunst, R. N., Lohmueller,
1233 K. E., Ralph, P. L., Schrider, D. R., Siepel, A., Kelleher, J., and Kern,
1234 A. D. (2020). A community-maintained standard library of population
1235 genetic models. *eLife*, 9:e54967.
- 1236 Alonso-Blanco, C., Andrade, J., Becker, C., Bemm, F., Bergelson, J., Borg-
1237 wardt, K. M., Cao, J., Chae, E., Dezwaan, T. M., Ding, W., et al. (2016).
1238 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis*
1239 *thaliana*. *Cell*, 166(2):481–491.
- 1240 Avise, J. C. (2009). Phylogeography: retrospect and prospect. *Journal of*
1241 *Biogeography*, 36(1):3–15.
- 1242 Barton, N. H. and Wilson, I. (1995). Genealogies and geography. *Philosophical*
1243 *Transactions of the Royal Society of London B*, 349(1327):49–59.
- 1244 Battey, C. J., Ralph, P. L., and Kern, A. D. (2020). Predicting geographic
1245 location from genetic variation with deep neural networks. *eLife*, 9:e54507.
- 1246 Bedford, T., Cobey, S., Beerli, P., and Pascual, M. (2010). Global migration
1247 dynamics underlie evolution and persistence of human influenza A (H3N2).
1248 *PLoS Pathogens*, 6(5):e1000918.
- 1249 Biek, R., Henderson, J. C., Waller, L. A., Rupprecht, C. E., and Real, L. A.
1250 (2007). A high-resolution genetic signature of demographic and spatial
1251 expansion in epizootic rabies virus. *Proceedings of the National Academy*
1252 *of Sciences*, 104(19):7993–7998.
- 1253 Bomblies, K., Yant, L., Laitinen, R. A., Kim, S.-T., Hollister, J. D., Warth-
1254 mann, N., Fitz, J., and Weigel, D. (2010). Local-scale patterns of genetic
1255 variability, outcrossing, and spatial structure in natural stands of *Ara-*
1256 *bidopsis thaliana*. *PLoS Genetics*, 6(3):e1000890.
- 1257 Bradburd, G. S. and Ralph, P. L. (2019). Spatial population genetics:
1258 it's about time. *Annual Review of Ecology, Evolution, and Systematics*,
1259 50:427–449.
- 1260 Cavalli-Sforza, L. L., Barrai, I., and Edwards, A. W. (1964). Analysis of
1261 human evolution under random genetic drift. In *Cold Spring Harbor sym-*
1262 *posia on quantitative biology*, volume 29, pages 9–20. Cold Spring Harbor
1263 Laboratory Press.

- 1264 Cavalli-Sforza, L. L. and Edwards, A. W. (1967). Phylogenetic analysis.
1265 models and estimation procedures. *American Journal of Human Genetics*,
1266 19:233–257.
- 1267 Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., and
1268 Lee, J. J. (2015). Second-generation PLINK: rising to the challenge of
1269 larger and richer datasets. *Gigascience*, 4(7):s13742–015–0047–8.
- 1270 Coop, G. (2013). How many genetic ancestors
1271 do I have? [https://gcbias.org/2013/11/11/
1272 how-does-your-number-of-genetic-ancestors-grow-back-over-time/](https://gcbias.org/2013/11/11/how-does-your-number-of-genetic-ancestors-grow-back-over-time/).
1273 [Online; accessed 1-July-2021].
- 1274 Coop, G. (2017). Where did your genetic ancestors come from? <https://gcbias.org/2017/12/19/1628/>. [Online; accessed 1-July-2021].
- 1275
- 1276 Dellicour, S., Durkin, K., Hong, S. L., Vanmechelen, B., Martí-Carreras, J.,
1277 Gill, M. S., Meex, C., Bontems, S., André, E., Gilbert, M., et al. (2021a). A
1278 phylodynamic workflow to rapidly gain insights into the dispersal history
1279 and dynamics of SARS-CoV-2 lineages. *Molecular Biology and Evolution*,
1280 38(4):1608–1613.
- 1281 Dellicour, S., Hong, S. L., Vrancken, B., Chaillon, A., Gill, M. S., Maurano,
1282 M. T., Ramaswami, S., Zappile, P., Marier, C., Harkins, G. W., et al.
1283 (2021b). Dispersal dynamics of SARS-CoV-2 lineages during the first epi-
1284 demic wave in New York City. *PLoS Pathogens*, 17(5):e1009571.
- 1285 Donnelly, K. P. (1983). The probability that related individuals share some
1286 section of genome identical by descent. *Theoretical Population Biology*,
1287 23(1):34–63.
- 1288 Durvasula, A., Fulgione, A., Gutaker, R. M., Alacakaptan, S. I., Flood, P. J.,
1289 Neto, C., Tsuchimatsu, T., Burbano, H. A., Picó, F. X., Alonso-Blanco, C.,
1290 et al. (2017). African genomes illuminate the early history and transition
1291 to selfing in *Arabidopsis thaliana*. *Proceedings of the National Academy of
1292 Sciences*, 114(20):5213–5218.
- 1293 Edwards, A. W. (1970). Estimation of the branch points of a branching
1294 diffusion process. *Journal of the Royal Statistical Society B*, 32(2):155–
1295 164.

- 1296 Exposito-Alonso, M., Becker, C., Schuenemann, V. J., Reiter, E., Setzer, C.,
1297 Slovak, R., Brachi, B., Hagmann, J., Grimm, D. G., Chen, J., et al. (2018).
1298 The rate and potential relevance of new mutations in a colonizing plant
1299 lineage. *PLoS Genetics*, 14(2):e1007155.
- 1300 Felsenstein, J. (1973). Maximum-likelihood estimation of evolutionary
1301 trees from continuous characters. *American Journal of Human Genetics*,
1302 25(5):471.
- 1303 Felsenstein, J. (1985). Phylogenies and the comparative method. *The Amer-
1304 ian Naturalist*, 125(1):1–15.
- 1305 Fulgione, A. and Hancock, A. M. (2018). Archaic lineages broaden our view
1306 on the history of *Arabidopsis thaliana*. *New Phytologist*, 219(4):1194–1198.
- 1307 Fulgione, A., Koornneef, M., Roux, F., Hermisson, J., and Hancock, A. M.
1308 (2018). Madeiran *Arabidopsis thaliana* reveals ancient long-range coloniza-
1309 tion and clarifies demography in Eurasia. *Molecular Biology and Evolution*,
1310 35(3):564–574.
- 1311 Grafen, A. (1989). The phylogenetic regression. *Philosophical Transactions
1312 of the Royal Society of London B*, 326(1233):119–157.
- 1313 Griffiths, R. C. and Tavare, S. (1994). Sampling theory for neutral alleles in
1314 a varying environment. *Philosophical Transactions of the Royal Society of
1315 London B*, 344(1310):403–410.
- 1316 Griffiths, R. C. and Tavare, S. (1997). Computational methods for the coa-
1317 lescent. *IMA Volumes in Mathematics and its Applications*, 87:165–182.
- 1318 Haller, B. C., Galloway, J., Kelleher, J., Messer, P. W., and Ralph,
1319 P. L. (2019). Tree-sequence recording in SLiM opens new horizons for
1320 forward-time simulation of whole genomes. *Molecular Ecology Resources*,
1321 19(2):552–566.
- 1322 Haller, B. C. and Messer, P. W. (2019). SLiM 3: Forward genetic simula-
1323 tions beyond the Wright–Fisher model. *Molecular Biology and Evolution*,
1324 36(3):632–637.
- 1325 Harmon, L. J. (2019). *Phylogenetic comparative methods*.
1326 <https://lukejharmon.github.io/pcm/>.

- 1327 Harmon, L. J., Losos, J. B., Jonathan Davies, T., Gillespie, R. G., Gittleman,
1328 J. L., Bryan Jennings, W., Kozak, K. H., McPeek, M. A., Moreno-Roark,
1329 F., Near, T. J., et al. (2010). Early bursts of body size and shape evolution
1330 are rare in comparative data. *Evolution*, 64(8):2385–2396.
- 1331 Hewitt, G. (2000). The genetic legacy of the quaternary ice ages. *Nature*,
1332 405(6789):907–913.
- 1333 Hsu, C.-W., Lo, C.-Y., and Lee, C.-R. (2019). On the postglacial spread of
1334 human commensal *Arabidopsis thaliana*: journey to the east. *New Phytologist*,
1335 222(3):1447–1457.
- 1336 Hudson, R. R. (2001). Two-locus sampling distributions and their applica-
1337 tion. *Genetics*, 159(4):1805–1817.
- 1338 Hudson, R. R. (2002). Generating samples under a Wright–Fisher neutral
1339 model of genetic variation. *Bioinformatics*, 18(2):337–338.
- 1340 Keightley, P. D. and Jackson, B. C. (2018). Inferring the probability of
1341 the derived vs. the ancestral allelic state at a polymorphic site. *Genetics*,
1342 209(3):897–906.
- 1343 Kelleher, J., Etheridge, A., Vber, A., and Barton, N. (2016a). Spread of
1344 pedigree versus genetic ancestry in spatially distributed populations. *The-
1345 oretical Population Biology*, 108:1–12.
- 1346 Kelleher, J., Etheridge, A. M., and McVean, G. (2016b). Efficient coalescent
1347 simulation and genealogical analysis for large sample sizes. *PLoS Compu-
1348 tational Biology*, 12(5):e1004842.
- 1349 Kelleher, J., Thornton, K. R., Ashander, J., and Ralph, P. L. (2018). Ef-
1350 ficient pedigree recording for fast population genetics simulation. *PLoS
1351 Computational Biology*, 14(11):e1006581.
- 1352 Kelleher, J., Wong, Y., Wohns, A. W., Fadil, C., Albers, P. K., and McVean,
1353 G. (2019). Inferring whole-genome histories in large population datasets.
1354 *Nature Genetics*, 51(9):1330–1338.
- 1355 Kim, S., Plagnol, V., Hu, T. T., Toomajian, C., Clark, R. M., Ossowski,
1356 S., Ecker, J. R., Weigel, D., and Nordborg, M. (2007). Recombina-
1357 tion and linkage disequilibrium in *Arabidopsis thaliana*. *Nature Genetics*,
1358 39(9):1151–1155.

- 1359 Kimura, M. (1980). A simple method for estimating evolutionary rates of
1360 base substitutions through comparative studies of nucleotide sequences.
1361 *Journal of Molecular Evolution*, 16(2):111–120.
- 1362 Knowles, L. L. (2009). Statistical phylogeography. *Annual Review of Ecology, Evolution, and Systematics*, 40:593–612.
- 1364 Lachance, J. (2009). Inbreeding, pedigree size, and the most recent common
1365 ancestor of humanity. *Journal of Theoretical Biology*, 261(2):238–247.
- 1366 Landis, M. J., Schraiber, J. G., and Liang, M. (2013). Phylogenetic analysis
1367 using Lévy processes: finding jumps in the evolution of continuous traits.
1368 *Systematic Biology*, 62(2):193–204.
- 1369 Larriba, F. and Fearnhead, P. (2011). On composite likelihoods in statistical
1370 genetics. *Statistica Sinica*, 21(1):43–69.
- 1371 Lee, C.-R., Svardal, H., Farlow, A., Exposito-Alonso, M., Ding, W.,
1372 Novikova, P., Alonso-Blanco, C., Weigel, D., and Nordborg, M. (2017).
1373 On the post-glacial spread of human commensal *Arabidopsis thaliana*. *Nature Communications*, 8(1):1–12.
- 1375 Lemey, P., Hong, S. L., Hill, V., Baele, G., Poletto, C., Colizza, V.,
1376 OToole, Á., McCrone, J. T., Andersen, K. G., Worobey, M., et al.
1377 (2020). Accommodating individual travel history and unsampled diversity
1378 in Bayesian phylogeographic inference of SARS-CoV-2. *Nature Communications*, 11(1):1–14.
- 1380 Lemey, P., Rambaut, A., Drummond, A. J., and Suchard, M. A. (2009).
1381 Bayesian phylogeography finds its roots. *PLoS Computational Biology*,
1382 5(9):e1000520.
- 1383 Lemey, P., Rambaut, A., Welch, J. J., and Suchard, M. A. (2010). Phy-
1384 logeography takes a relaxed random walk in continuous space and time.
1385 *Molecular Biology and Evolution*, 27(8):1877–1885.
- 1386 Malécot, G. (1948). *Les mathématiques de l'hérédité*. Masson.
- 1387 Martin, M. A., VanInsberghe, D., and Koelle, K. (2021). Insights from SARS-
1388 CoV-2 sequences. *Science*, 371(6528):466–467.

- 1389 Meligkotsidou, L. and Fearnhead, P. (2007). Postprocessing of genealogical
1390 trees. *Genetics*, 177(1):347–358.
- 1391 Nordborg, M. (2000). Linkage disequilibrium, gene trees and selfing: an
1392 ancestral recombination graph with partial self-fertilization. *Genetics*,
1393 154(2):923–929.
- 1394 Novembre, J. and Slatkin, M. (2009). Likelihood-based inference in isolation-
1395 by-distance models using the spatial distribution of low-frequency alleles.
1396 *Evolution*, 63(11):2914–2925.
- 1397 O'Meara, B. C., Ané, C., Sanderson, M. J., and Wainwright, P. C. (2006).
1398 Testing for different rates of continuous trait evolution using likelihood.
1399 *Evolution*, 60(5):922–933.
- 1400 Platt, A., Horton, M., Huang, Y. S., Li, Y., Anastasio, A. E., Mulyati,
1401 N. W., Ågren, J., Bossdorf, O., Byers, D., Donohue, K., et al. (2010).
1402 The scale of population structure in *Arabidopsis thaliana*. *PLoS Genetics*,
1403 6(2):e1000843.
- 1404 Ralph, P. and Coop, G. (2013). The geography of recent genetic ancestry
1405 across europe. *PLoS Biology*, 11(5):e1001555.
- 1406 Ralph, P., Thornton, K., and Kelleher, J. (2020). Efficiently summarizing
1407 relationships in large samples: a general duality between statistics of ge-
1408 nealogies and genomes. *Genetics*, 215(3):779–797.
- 1409 Rannala, B. and Yang, Z. (1996). Probability distribution of molecular evo-
1410 lutionary trees: a new method of phylogenetic inference. *Journal of Molec-*
1411 *ular Evolution*, 43(3):304–311.
- 1412 Rasmussen, M. D., Hubisz, M. J., Gronau, I., and Siepel, A. (2014).
1413 Genome-wide inference of ancestral recombination graphs. *PLoS Genetics*,
1414 10(5):e1004342.
- 1415 Reich, D. (2018). *Who we are and how we got here: ancient DNA and the*
1416 *new science of the human past*. Oxford University Press.
- 1417 Ringbauer, H., Coop, G., and Barton, N. H. (2017). Inferring recent demog-
1418 raphy from isolation by distance of long shared sequence blocks. *Genetics*,
1419 205(3):1335–1351.

- 1420 Rousset, F. (1997). Genetic differentiation and estimation of gene flow from
1421 F-statistics under isolation by distance. *Genetics*, 145(4):1219–1228.
- 1422 Salomé, P., Bomblies, K., Fitz, J., Laitinen, R., Warthmann, N., Yant,
1423 L., and Weigel, D. (2012). The recombination landscape in *Arabidopsis*
1424 *thaliana* F2 populations. *Heredity*, 108(4):447–455.
- 1425 Schlüter, D., Price, T., Mooers, A. Ø., and Ludwig, D. (1997). Likelihood of
1426 ancestor states in adaptive radiation. *Evolution*, 51(6):1699–1711.
- 1427 Shirsekar, G., Devos, J., Latorre, S. M., Blaha, A., Dias, M. Q., Hernando,
1428 A. G., Lundberg, D. S., Burbano, H. A., Fenster, C. B., and Weigel,
1429 D. (2021). Fine-scale population structure of north american *Arabidop-*
1430 *sis thaliana* reveals multiple sources of introduction from across Eurasia.
1431 *bioRxiv*, doi:10.1101/2021.01.22.427575.
- 1432 Slatkin, M. and Racimo, F. (2016). Ancient DNA and human history. *Pro-*
1433 *ceedings of the National Academy of Sciences*, 113(23):6380–6387.
- 1434 Speidel, L., Cassidy, L., Davies, R. W., Hellenthal, G., Skoglund, P., and
1435 Myers, S. (2021). Inferring population histories for ancient genomes using
1436 genome-wide genealogies. *bioRxiv*, doi:10.1101/2021.02.17.431573.
- 1437 Speidel, L., Forest, M., Shi, S., and Myers, S. R. (2019). A method for
1438 genome-wide genealogy estimation for thousands of samples. *Nature Ge-*
1439 *netics*, 51(9):1321–1329.
- 1440 Stephens, M. and Donnelly, P. (2000). Inference in molecular population
1441 genetics. *Journal of the Royal Statistical Society B*, 62(4):605–635.
- 1442 Stern, A. J., Speidel, L., Zaitlen, N. A., and Nielsen, R. (2021). Disentan-
1443 gling selection on genetically correlated polygenic traits via whole-genome
1444 genealogies. *The American Journal of Human Genetics*, 108(2):219–239.
- 1445 Stern, A. J., Wilton, P. R., and Nielsen, R. (2019). An approximate full-
1446 likelihood method for inferring selection and allele frequency trajectories
1447 from DNA sequence data. *PLoS Genetics*, 15(9):e1008384.
- 1448 Sukumaran, J. and Holder, M. T. (2010). DendroPy: a Python library for
1449 phylogenetic computing. *Bioinformatics*, 26(12):1569–1571.

- 1450 Varin, C., Reid, N., and Firth, D. (2011). An overview of composite likelihood
1451 methods. *Statistica Sinica*, 21(1):5–42.
- 1452 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T.,
1453 Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van
1454 der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson,
1455 A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, I., Feng, Y.,
1456 Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R.,
1457 Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro,
1458 A. H., Pedregosa, F., van Mulbregt, P., and SciPy 1.0 Contributors (2020).
1459 SciPy 1.0: Fundamental algorithms for scientific computing in Python.
1460 *Nature Methods*, 17:261–272.
- 1461 Volz, E. M., Koelle, K., and Bedford, T. (2013). Viral phylodynamics. *PLoS
1462 Computational Biology*, 9(3):e1002947.
- 1463 Wakeley, J., King, L., Low, B. S., and Ramachandran, S. (2012). Gene
1464 genealogies within a fixed pedigree, and the robustness of Kingmans coa-
1465 lescent. *Genetics*, 190(4):1433–1445.
- 1466 Wang, I. J. and Bradburd, G. S. (2014). Isolation by environment. *Molecular
1467 Ecology*, 23(23):5649–5662.
- 1468 Wilkins, J. F. (2004). A separation-of-timescales approach to the coalescent
1469 in a continuous population. *Genetics*, 168(4):2227–2244.
- 1470 Wilkins, J. F. and Wakeley, J. (2002). The coalescent in a continuous, finite,
1471 linear population. *Genetics*, 161(2):873–888.
- 1472 Wiuf, C. (2006). Consistency of estimators of population scaled parameters
1473 using composite likelihood. *Journal of Mathematical Biology*, 53(5):821–
1474 841.
- 1475 Wohns, A. W., Wong, Y., Jeffery, B., Akbari, A., Mallick, S., Pin-
1476 hasi, R., Patterson, N., Reich, D., Kelleher, J., and McVean, G.
1477 (2021). A unified genealogy of modern and ancient genomes. *bioRxiv*,
1478 doi:10.1101/2021.02.16.431497.
- 1479 Worobey, M., Pekar, J., Larsen, B. B., Nelson, M. I., Hill, V., Joy, J. B.,
1480 Rambaut, A., Suchard, M. A., Wertheim, J. O., and Lemey, P. (2020).

- 1481 The emergence of SARS-CoV-2 in Europe and North America. *Science*,
1482 370(6516):564–570.
- 1483 Wright, S. (1943). Isolation by distance. *Genetics*, 28(2):114–138.
- 1484 Zeng, L., Gu, Z., Xu, M., Zhao, N., Zhu, W., Yonezawa, T., Liu, T., Qiong,
1485 L., Tersing, T., Xu, L., et al. (2017). Discovery of a high-altitude ecotype
1486 and ancient lineage of *Arabidopsis thaliana* from tibet. *Science Bulletin*,
1487 62(24):1628–1630.
- 1488 Zou, Y.-P., Hou, X.-H., Wu, Q., Chen, J.-F., Li, Z.-W., Han, T.-S., Niu, X.-
1489 M., Yang, L., Xu, Y.-C., Zhang, J., et al. (2017). Adaptation of *Arabidopsis*
1490 *thaliana* to the Yangtze River basin. *Genome Biology*, 18(1):1–11.

1491 **Supplementary Text**

1492 **Multi-epoch dispersal rate estimates**

1493 In both cases the two-epoch maximum negative log-likelihood was much
1494 larger than the maximum negative log-likelihood from the one-epoch model
1495 (mean differences of > 500 in B and > 2000 in C with true trees and > 700
1496 and > 400 with inferred trees), providing strong statistical support for the
1497 two-epoch model. However, when we fit two-epoch models to simulations
1498 with a constant dispersal rate (Figure S3A,D), the two-epoch models still
1499 tended to have larger negative log-likelihoods than the one-epoch models
1500 (mean differences of ≈ 500 for $\sigma_x = \sigma_y = 0.25$ and ≈ 50 for $\sigma_x = \sigma_y = 0.5$
1501 with true trees and ≈ 200 and ≈ 250 with inferred trees). Using the true
1502 trees, at higher dispersal rates ($\sigma_x = \sigma_y = 0.5$, Figure S3D) the estimates
1503 of dispersal in each of the two epochs were similar to one another, and to
1504 the one-epoch estimate. This likely accounts for the relatively small in-
1505 crease in likelihood under the two-epoch model. At lower dispersal rates
1506 ($\sigma_x = \sigma_y = 0.25$, Figure S3A) we see a different pattern, where the true (and
1507 inferred) trees now tend to overestimate dispersal in the more distant epoch.
1508 We suspect this to be a sampling bias; in cutting trees off at 10^3 generations
1509 we effectively give more weight to lineages that coalesce quickly, creating an
1510 upwards bias in dispersal rates. Supporting this hypothesis, using a deeper
1511 cutoff of 10^4 generations produced more accurate estimates at these lower dis-
1512 persal rates (results not shown). Curiously, the inferred trees tell a different
1513 two-epoch story than the true trees under high dispersal (Figure S3D), but
1514 not under low dispersal (Figure S3A). This suggests tree inference is worse,
1515 or that errors have more impact on dispersal estimates, when the samples
1516 are more closely related.

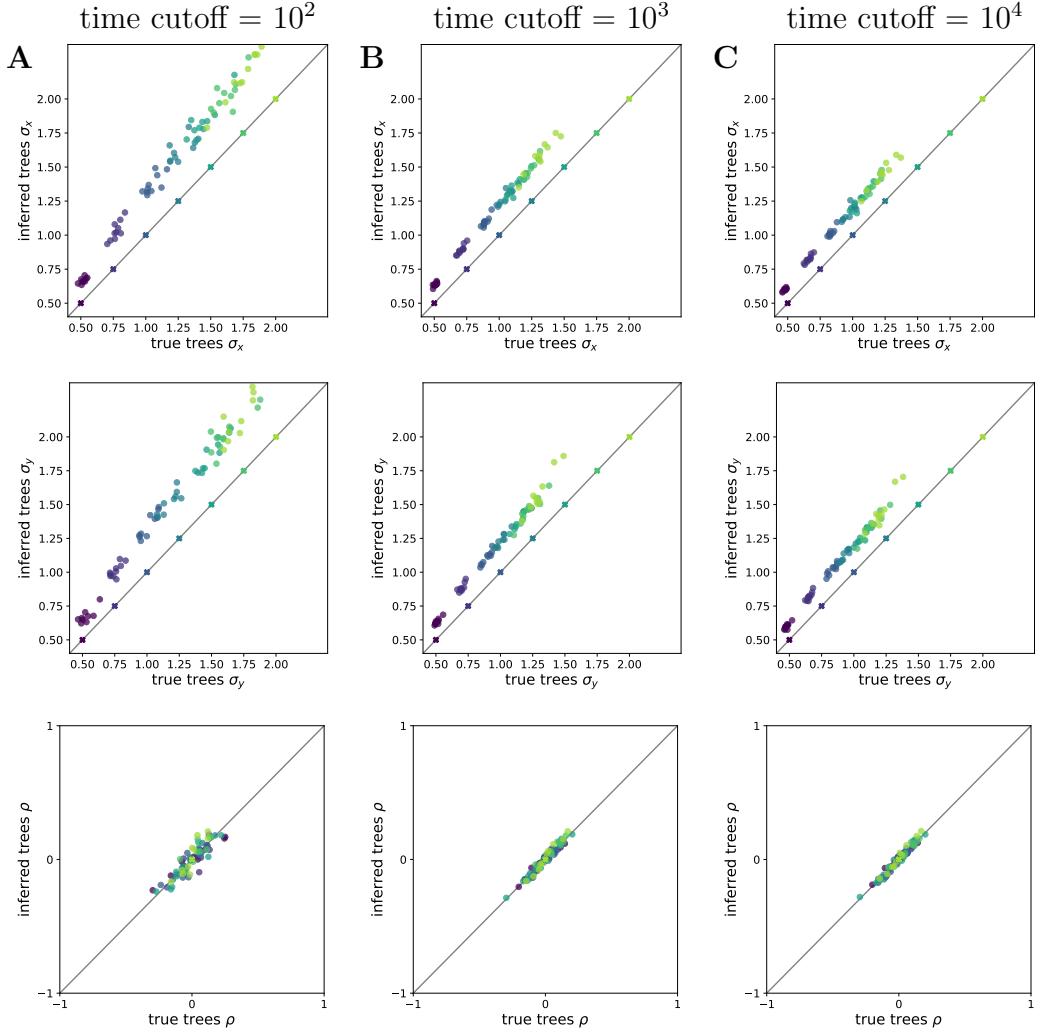


Figure S1: The effect of the time cutoff on dispersal estimates in one-epoch simulations. Panel **A** is as in the main text (Figure 2A), where we ignore everything beyond 10^2 generations ago (and here we show dispersal along the y-axis in the correlation in x and y). Panels **B** and **C** increase the time cutoff by one and two orders of magnitude, respectively. Dispersal estimates shrink as the time cutoff grows because smaller dispersal rates make it more likely all lineages have remained within the finite habitat (with reflecting boundaries). The variance in estimates seems to decline with the time cutoff, as larger cutoffs mean we use more information from the trees.

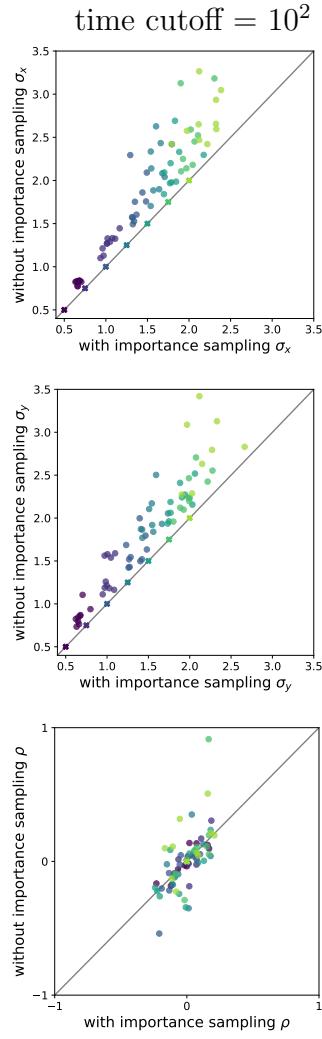


Figure S2: **The effect of importance sampling on dispersal estimates in one-epoch simulations.** In all cases importance sampling leads to smaller overestimates of the simulated dispersal rate (σ_x and σ_y) and generally exhibits lower variance in estimates across replicates (especially at high simulated dispersal rates). The ‘without importance sampling’ estimate comes from using only the first sample of branch lengths, while the ‘with importance sampling’ estimate uses all 10 samples.

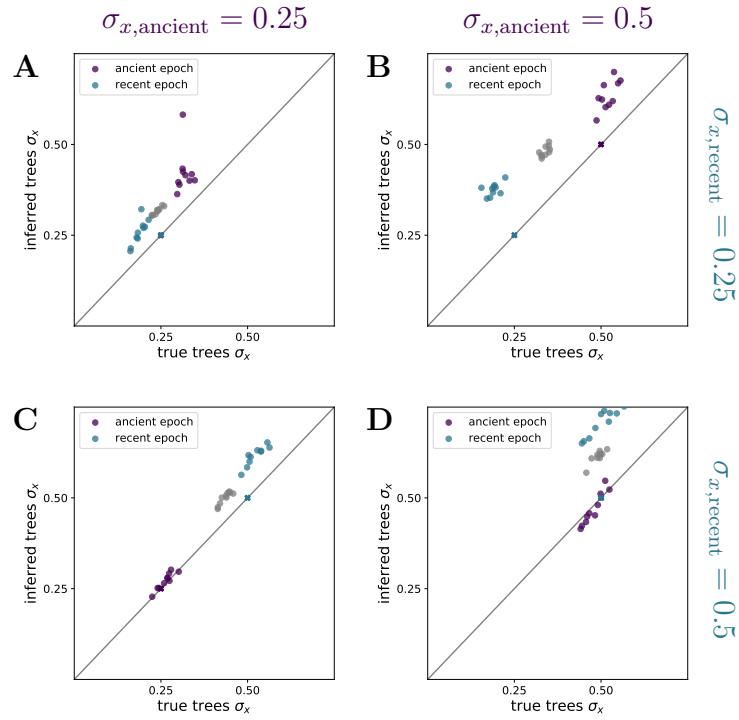


Figure S3: **Estimating two-epoch dispersal rates in simulations.** Panels **B** and **D** are panels **B** and **C**, respectively, in Figure 2. Panels **A** and **D** are simulated with a constant dispersal rate, $\sigma = 0.25$ and 0.5 , respectively. The coloured dots show the estimates for the recent and ancient epochs. The grey dots show the estimates under a one-epoch model.

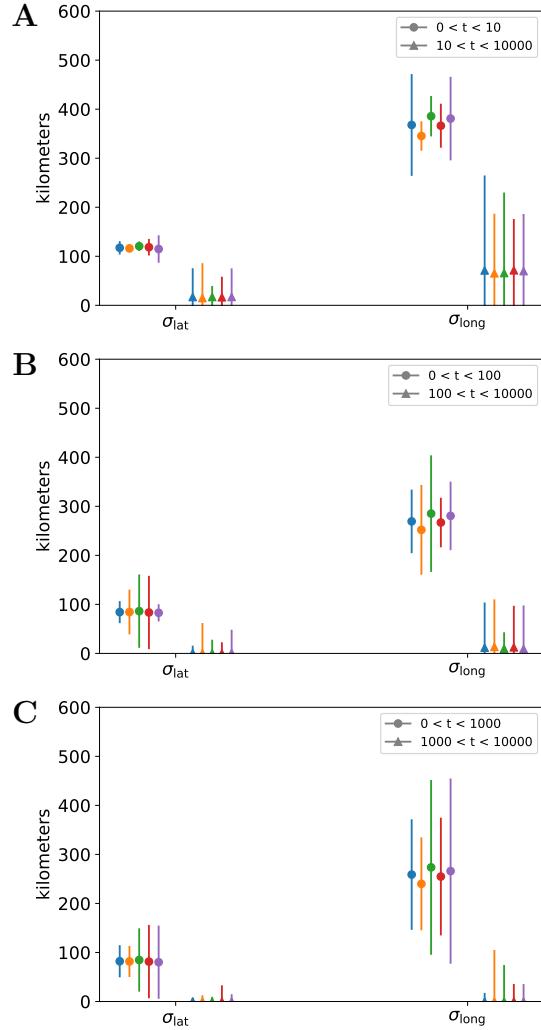


Figure S4: Dispersal rate estimates in *Arabidopsis thaliana* under a two-epoch model. Like Figure 3B, but here we also show the (less likely) split times of (B) 10^2 and (C) 10^3 .

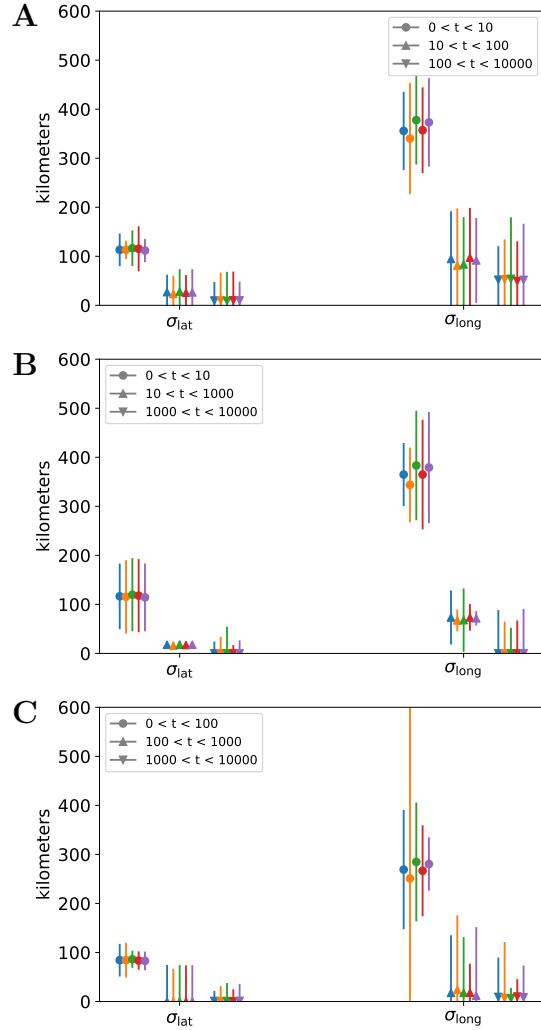
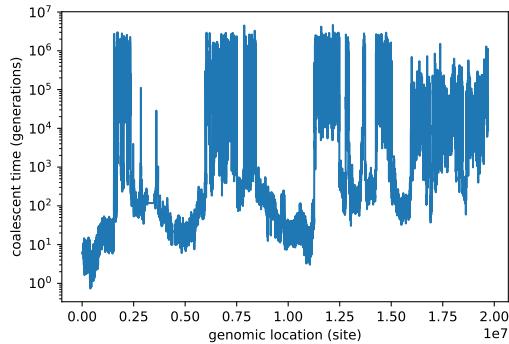
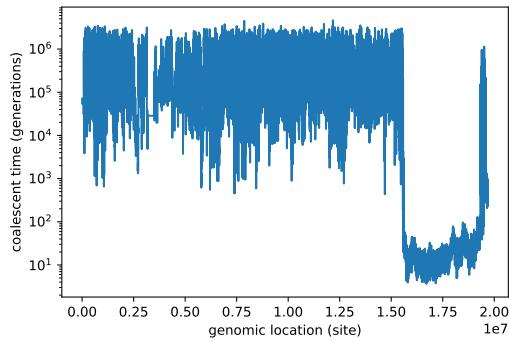


Figure S5: Dispersal rate estimates in *Arabidopsis thaliana* under a three-epoch model. Like Figure 3B, but here we also show the (less likely) split times of (A) $[10, 10^2]$ and (C) $[10^2, 10^3]$.

A) accessions 9737 & 9627 (chromosome 2)



B) accessions 9933 & 10015 (chromosome 2)



C) accessions 7314 & 6981 (chromosome 5)

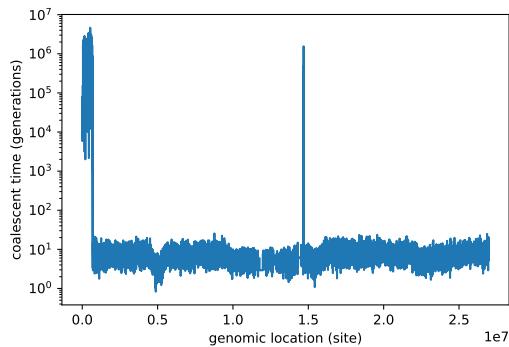


Figure S6: Coalescence times between particular pairs of samples along a particular chromosome. Calculated using tskit's `divergence()` function.

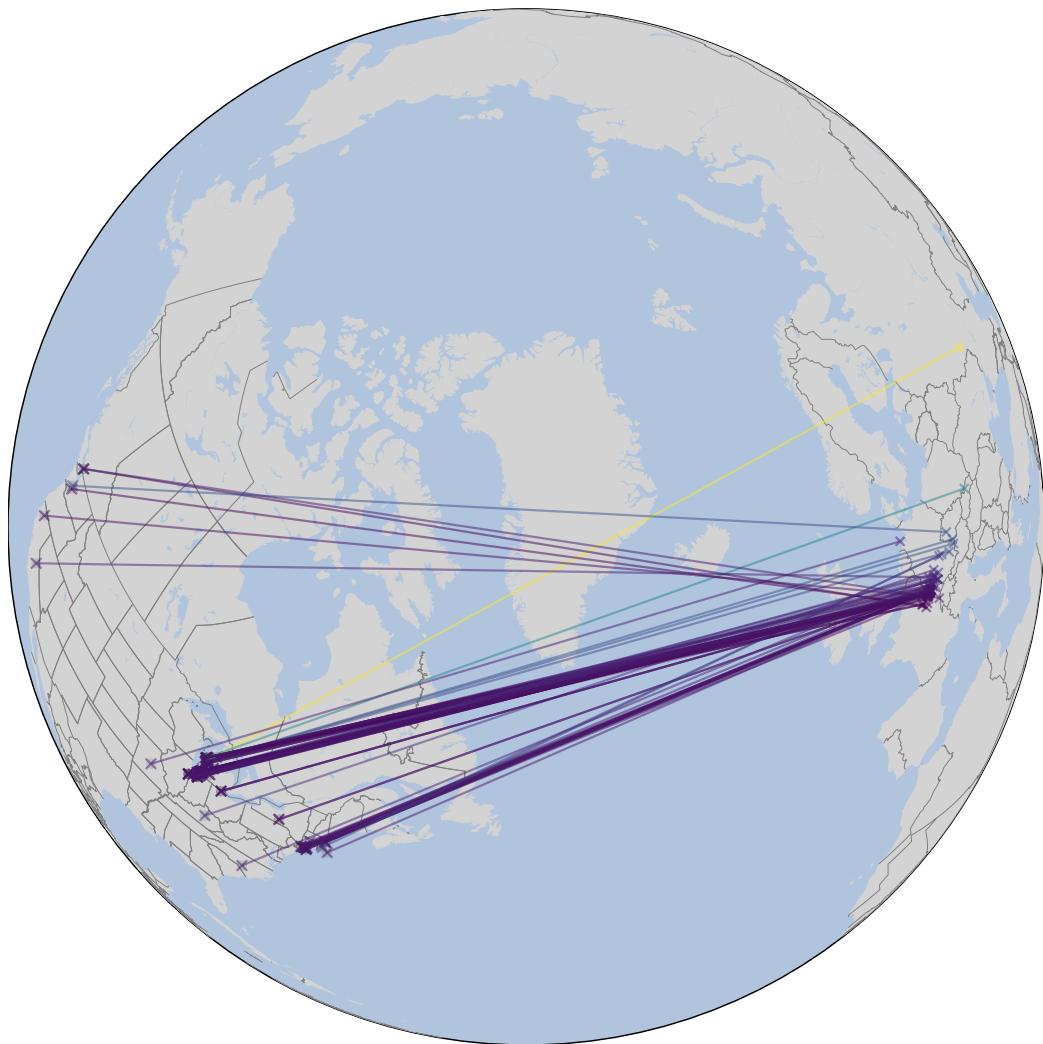
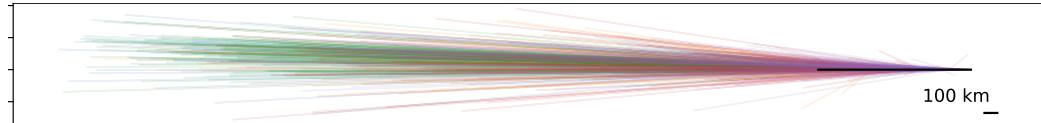
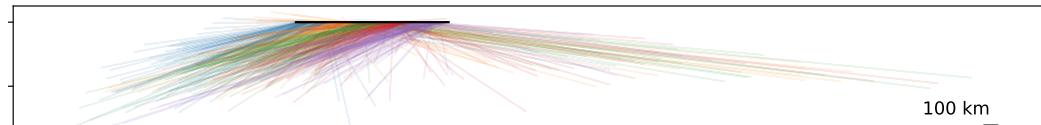


Figure S7: Connecting the sample locations with the inferred present-day locations of the North American samples. Color indicates longitude of the inferred location. See Figure 6 for more details.

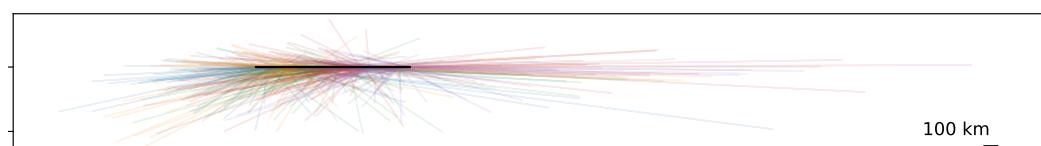
accession 6830



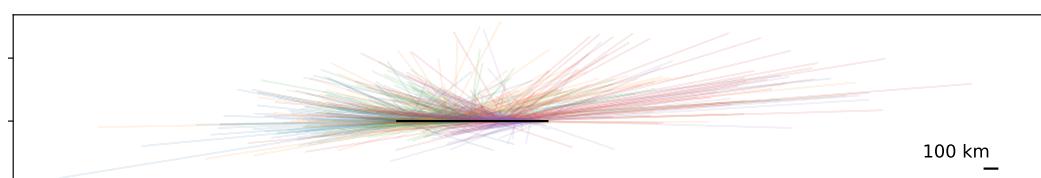
accession 7126



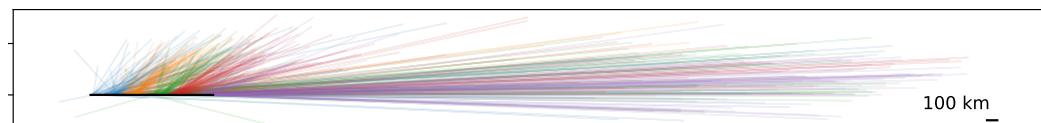
accession 7394



accession 9743



accession 9933



accession 9993

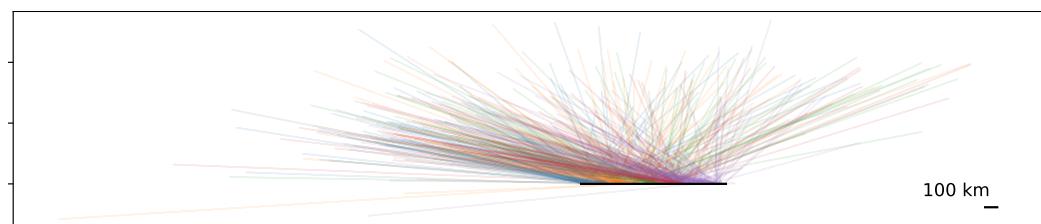


Figure S8: **Genome-view of ancestral displacements for all 6 samples in Figure 9.** See Figure 9 for more details.