

Predicting Geographic Location from Genetic Variation with Deep Neural Networks

C.J. Battey¹, Peter L. Ralph¹, and Andrew D. Kern¹

¹University of Oregon Institute of Ecology and Evolution

December 11, 2019

Abstract

Most organisms are more closely related to nearby than distant members of their species, creating spatial autocorrelations in genetic data. This allows us to predict the location of origin of a genetic sample by comparing it to a set of samples of known geographic origin. Here we describe a deep learning method, which we call *Locator*, to accomplish this task faster and more accurately than existing approaches. In simulations, *Locator* infers sample location to within 4.1 generations of dispersal and runs at least an order of magnitude faster than a recent model-based approach. We leverage *Locator*'s computational efficiency to predict locations separately in windows across the genome, which allows us to both quantify uncertainty and describe the mosaic ancestry and patterns of geographic mixing that characterize many populations. Applied to whole-genome sequence data from *Plasmodium* parasites, *Anopheles* mosquitoes, and global human populations, this approach yields median test errors of 16.9km, 5.7km, and 85km, respectively.

Introduction

In natural populations, local mate selection and dispersal create correlations between geographic location and genetic variation – each individual's genome is a mosaic of material inherited from recent ancestors that are usually geographically nearby. Given a set of genotyped individuals of known geographic provenance, it is therefore possible to predict the location of new samples from genetic information alone (Guillot et al., 2015; Yang et al., 2012; Wasser et al., 2004; Rañola et al., 2014; Bhaskar et al., 2016; Baran et al., 2013). This task has forensic applications – for example, estimating the location of trafficked elephant ivory as in Wasser et al. (2004) – and also offers a way to analyze variation in geographic ancestry without assuming the existence of discrete ancestral populations.

The most common approaches to estimating sample locations are based on unsupervised genotype clustering or dimensionality reduction techniques. Genetic data from samples of both known and unknown origin are jointly analyzed, and unknown samples are assigned to the location of known individuals with which they share a genotype cluster or region of PC space (Breidenbach et al., 2019; Battey et al., 2018; Cong et al., 2019). However, these methods require an additional mapping from genotype clusters or PC space to geography, and can produce nonsensical results if unknown samples are hybrids or do not originate from any of the sampled reference populations.

Existing methods for estimating sample location that explicitly model continuous landscapes use a two-step procedure. A smoothed map describing variation in allele frequencies over space

37 is first estimated for each allele based on the genotypes of individuals with known locations, and
38 locations of new samples are then predicted by maximizing the likelihood of observing a given
39 combination of alleles at the predicted location. In methods like SPASIBA (Guillot et al., 2015)
40 and SCAT (Wasser et al., 2004), allele frequency surfaces are estimated by fitting parameters of
41 a Gaussian function of set form (but see Rañola et al. (2014) for an alternate approach based on
42 smoothing techniques from image analysis).

43 Since all such methods use relatedness to other contemporary samples, any information about
44 the location of a new sample necessarily comes from ancestors shared with the reference set. As
45 illustrated in Figure 1, we expect *a priori*, that the genealogical relationships among a set of
46 samples (and therefore the spatial location of ancestors) will vary along the genome. This means
47 that a complete look at geographic ancestry would include not just a point estimate of spatial
48 location, but an estimate of uncertainty that accounts for the partially correlated genealogies of
49 recombining chromosomes.

50 In the past few years there has been an explosion in the use of supervised machine learning
51 for population genetics for a number of tasks, including detecting selection (Schridder and Kern,
52 2016; Mughal and DeGiorgio, 2018; Sugden et al., 2018), inferring admixture (Schridder et al., 2018;
53 Durvasula and Sankararaman, 2019), and performing demographic model selection (Pudlo et al.,
54 2015; Villanea and Schraiber, 2019). Applications to population genetics increasingly make use
55 of the latest generation of machine learning tools: deep neural networks (a.k.a. “deep learning”)
56 (Sheehan and Song, 2016; Kern and Schridder, 2018; Chan et al., 2018; Fligel et al., 2018; Adrion
57 et al., 2019). A significant feature of neural networks is that they allow the input of raw genotype
58 information, as we perform below, without initial compression into summary statistics.

59 In this paper, we introduce **Locator**, a highly efficient deep learning method for the prediction of
60 geographic origin of individuals from unphased genotype data. **Locator** uses deep neural networks
61 to perform prediction directly from genotypes, but without assuming any explicit model of how
62 genotypes vary over the landscape. Moreover, unlike many modern supervised machine learning
63 methods in population genetics, (e.g., Kern and Schridder, 2018) our training set need not be
64 obtained via simulation. We assume only that there is some function relating geographic locations
65 to the probability of observing a given combination of alleles, and use a deep, fully-connected neural
66 network to approximate this mapping for a set of genotyped individuals with known locations. The
67 trained network is then evaluated against a set of known individuals held out from the training
68 routine and used to predict the geographic location of new samples based on their genotypes.
69 Applied separately to windows across the genome, **Locator** also estimates uncertainty in individual-
70 level predictions, and can reveal portions of an individual’s genome enriched for ancestry from
71 specific geographic areas.

72 For the empirical population genomic data we analyze here, **Locator** achieves state-of-the-art
73 accuracy an order of magnitude faster than competing methods. Here we describe the implemen-
74 tation, test on simulated data, and demonstrate its use in empirical data by estimating sampling
75 locations for *Anopheles* mosquitoes in Africa from the AG1000G project (The Anopheles gambiae
76 1000 Genomes Consortium, 2015), *P. falciparum* parasites from Asia, Africa, and the Americas
77 from the *P. falciparum* community project (Pearson et al., 2019), and global human populations
78 from the Human Genome Diversity Project (HGDP; Bergström et al. (2019)).

79 Results

80 **Locator is fast and accurate**

81 We first evaluated **Locator**’s performance in simulations of populations evolving in continuous
82 space with varying rates of dispersal – an idealized setting in which all alleles should vary smoothly
83 over the map. In Figure 2 we show that validation error increases along with the dispersal rate

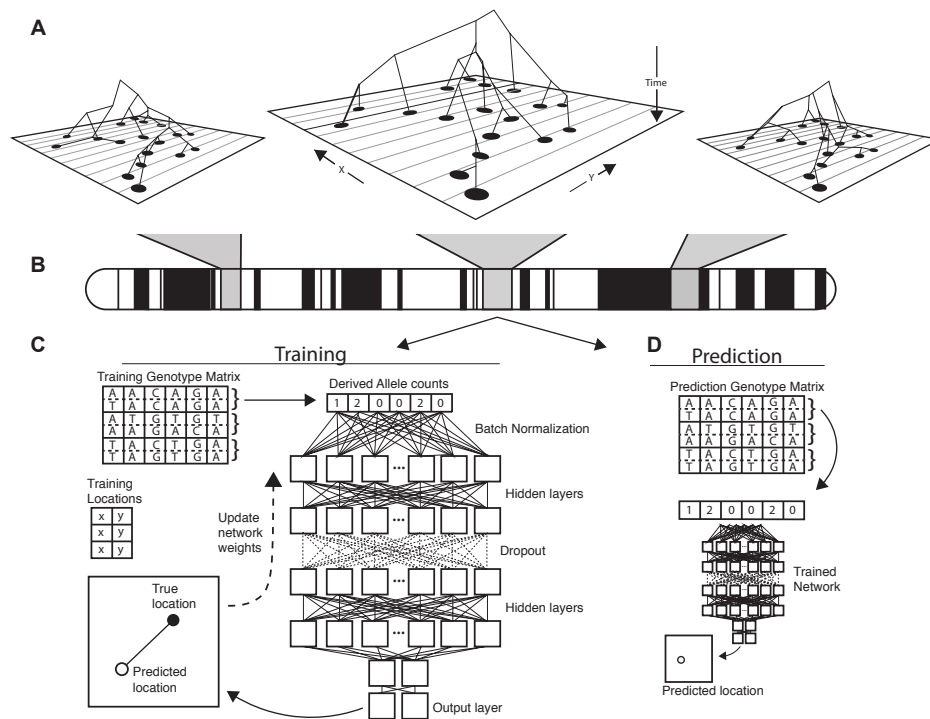


Figure 1: Conceptual schematic of our approach. Regions of the genome reflect correlated sets of genealogical relationships (A), each of which represents a set of ancestors with varying spatial positions back in time. We extract genotypes from windows across the genome (B), and train a deep neural network to approximate the relationship between genotypes and locations using Euclidean distance as the loss function (C). We can then use the trained network to predict the location of new genotypes held out from the training routine (D).

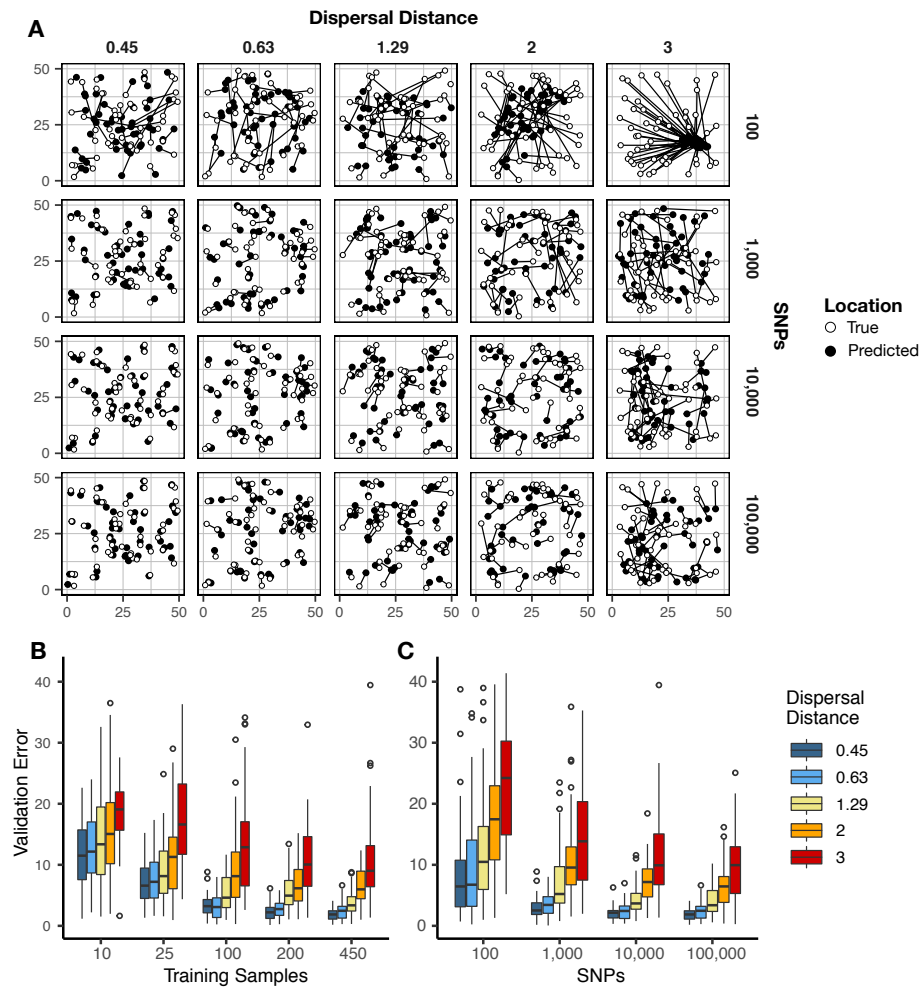


Figure 2: Validation error for *Locator* runs on simulations with varying neighborhood size. Simulations were on a 50 x 50 landscape and error is expressed in map units. A: True and predicted locations by neighborhood size and number of SNPs. 450 randomly-sampled individuals were used for training. B: Error for runs with 100,000 SNPs and varying numbers of training samples. C: Error for runs with 450 training samples and varying number of SNPs. Plots with error in terms of generations of expected dispersal are shown in Figure S2.

84 of the population. Interestingly, error is roughly constant when correcting for the dispersal rate
85 in each simulation, ranging from 3.16 to 4.09 generations of dispersal given our largest training
86 dataset (450 samples, 100,000 SNPs; Figure S2). This suggests that error primarily reflects the
87 underlying biological processes of dispersal and mate selection rather than simple noise from model
88 fitting.

89 Increasing the number of training samples or the number of SNPs improves accuracy for all
90 simulations (Figure 2B). However, we observed diminishing returns when using over 10,000 SNPs
91 or over 200 training samples. Median error for all simulations was also below 10 generations of
92 dispersal for all but the least-dispersive simulation using just 25 training samples; suggesting that
93 even relatively small training datasets can allow inference of broad-scale spatial locations. We
94 discuss theoretical limits on the accuracy of genetic location estimation in Appendix 1 .

95 We were interested to compare the performance of **Locator** to that of **SPASIBA** (Guillot et al.,
96 2015), the current state-of-the-art method for geographic prediction of genotype data (Figure 3).
97 However, we were unable to successfully run **SPASIBA** with 100,000 or more SNPs from a simulated
98 dataset or on simulations with dispersal rates of 0.63 or 1.29 map units/generation, due to out-
99 of-memory errors on a 64-bit system with 400Gb of RAM. We could however compare at smaller
100 numbers of SNPs and reduced dispersal. At a mean dispersal distance of 0.45 map units **SPASIBA**'s
101 median test error was slightly lower when run on 1,000 SNPs (Wilcoxon test, $p=0.009$) but results
102 were similar at 100 or 10,000 SNPs. (Wilcoxon test, $p = 0.184$ and 0.936). However, **Locator** is
103 much faster – training on 10,000 SNPs in less than two minutes while **SPASIBA** requires around
104 six and a half hours (Figure 2). These long run times are caused in part by the large number of
105 training localities in our simulated data, because **SPASIBA**'s run time scales with the product of
106 the number of genetic variants and the number of training localities (Guillot et al., 2015).

107 While the simulations conform well to modeling assumptions of most methods, we can also
108 compare performance on empirical data. By way of example, we applied **Locator** and **SPASIBA** to
109 subsets of SNPs from the first five million base pairs of chromosome 2L from the Ag1000G dataset
110 Miles and Harding (2017) (figure 3). **Locator** achieves much lower mean error on all runs with
111 more than 100 SNPs, and runs from 3.1x to 532x faster, depending on the number of SNPs. Maps
112 of predictions from both methods are shown in Figure S5. Extrapolating from these run times,
113 running a windowed whole-genome analysis of *Anopheles* in **SPASIBA** would require roughly 70
114 days of computation on an 80-CPU system for model training alone, versus 3.2 hours on one GPU
115 for **Locator**.

116 **Uncertainty and Variation along the Genome**

117 By running **Locator** in windows across the genome we aim to integrate over error associated with
118 the model training procedure while also representing the inherent uncertainty caused by spatial
119 drift of ancestral lineages backwards in time (Kelleher et al., 2016). This produces a cloud of
120 predicted locations distributed around the true sample location (Figure 4). For individuals near
121 the center of the landscape these clouds are roughly symmetrical, as expected from our model.
122 Predictions for individuals close to the edge of the landscape appear slightly asymmetrical and
123 are bounded by the true landscape edges, suggesting that our networks have learned the rough
124 shape of the sampled range. The true location was within the 50% contour of a 2d-kernel density
125 surface estimated from the set of per-window predictions for all test samples, demonstrating that
126 this distribution is indeed centered on the true location. We also tested the alternate approach of
127 bootstrapping over a single set of SNPs, which could be useful for smaller datasets or those lacking
128 a reference alignment. Results for this method are discussed in Supplementary figure S4.

129 Windowed analyses for the three empirical systems we studied are shown in the bottom panels
130 of Figures 5–7. We discuss the implications of these predictions for each species below, but in
131 general we find that the windowed analysis accurately describes uncertainty in a sample's location

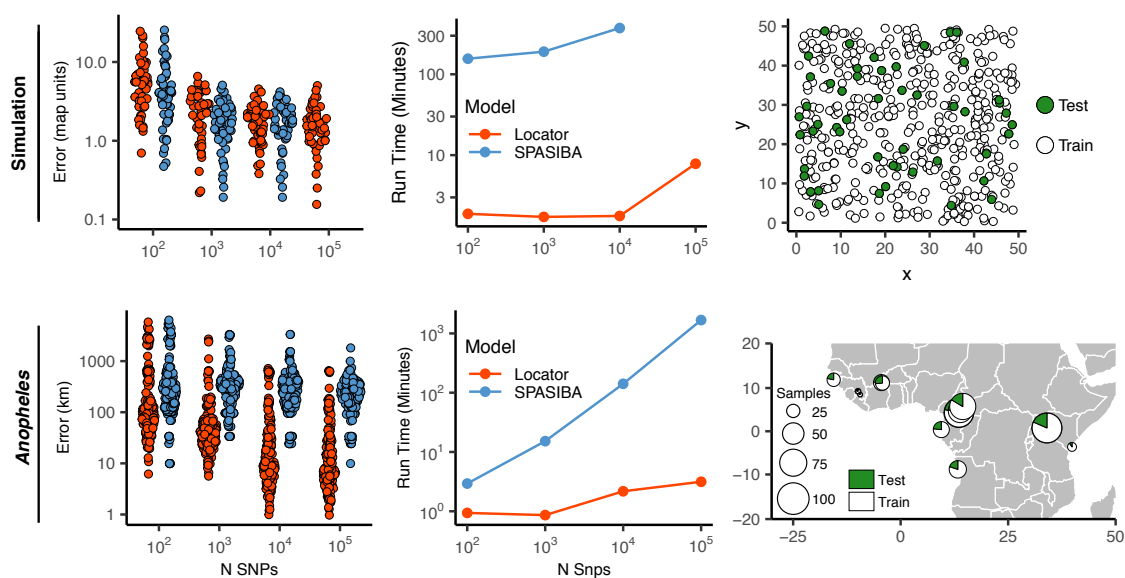


Figure 3: Test error and run times for *Locator* and *SPASIBA* on simulated data with dispersal distance equal to 0.45 map units/generation (top; 450 randomly sampled training samples) and empirical data from the ag1000g phase 1 dataset (bottom; 612 training samples from 14 sampling localities).

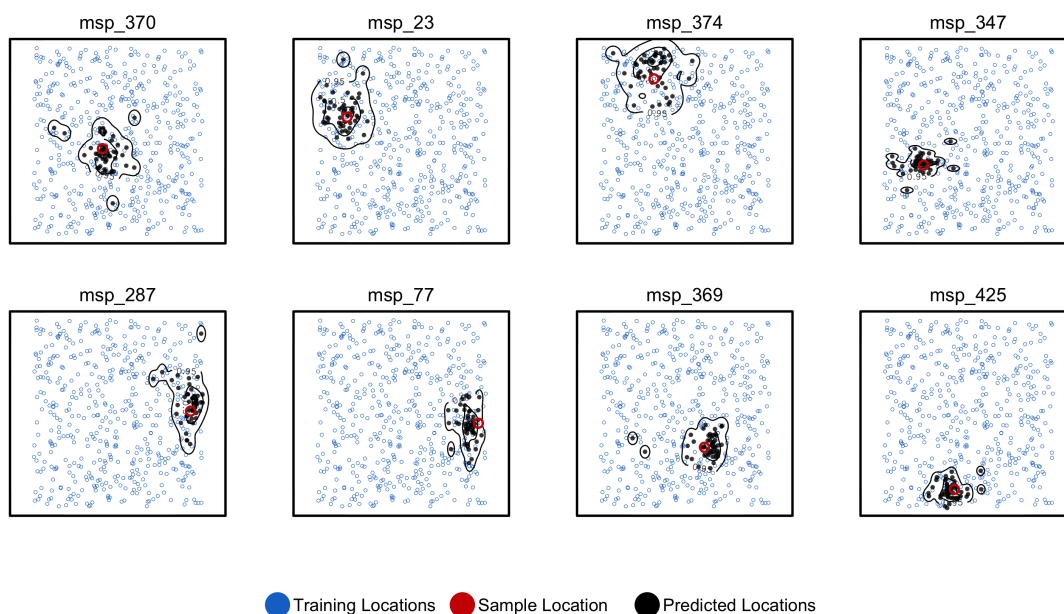


Figure 4: Predicted and true locations for 8 individuals simulated in a population with neighborhood size ≈ 25 . Black points are predictions from 2Mbp windows, blue points are training sample locations, and the red point is the true location for each individual. Contours show the 95%, 50%, and 10% quantiles of a two-dimensional kernel density across all windows.

132 – either surrounding a single location for samples with low error, or distributed across a wide
133 region including multiple training localities for samples with high error. In several cases predicted
134 locations also project in the direction of known historic migrations (as in human data), or are split
135 among localities shown in previous analyses to experience high gene flow (as in *Anopheles*).

136 We summarize genome-wide window predictions in two ways: 1) by taking a kernel density
137 estimate of the predictions and then finding the point of maximum density, and 2) by computing
138 the centroid of the windowed predictions. These estimates are similar in spirit to ensemble predic-
139 tion methods (Ho, 1995; Breiman, 1996). In general we found that the maximum kernel density
140 estimator has lower error, but tends to show classification behavior more than the centroid esti-
141 mator – snapping to a single training locality rather than interpolating between sets of localities
142 for samples with variable window predictions.

143 Empirical Analysis

144 *Anopheles* mosquitoes

145 We next turn our attention to the application of *Locator* to empirical population genomic datasets.
146 In Figure 5 we show predicted and true locations for 153 individuals from the Ag1000g dataset of
147 *Anopheles gambiae* and *A. coluzzii*, estimated in 2Mbp windows across the genome. The location
148 with highest kernel density across all windows had a median error of 5.7km, and the centroid of the
149 per-window predictions had a median error of 36 km (Table S2). Significant prediction error occurs
150 only between sites in Cameroon, Burkina Faso, and the Republic of Guinea – localities which were
151 also assigned to a single ancestry cluster in the ADMIXTURE analysis in Miles et al. (2017).
152 However uncertainty for these samples was relatively well described by visualizing the spread of
153 per-window predictions, with predicted locations generally lying between sets of localities. The
154 true locality was within the 95% interval of the kernel density across all windows for all samples.

155 *Plasmodium falciparum*

156 In a windowed analysis of *P. falciparum*, *Locator*'s median error is 16.92 km using the maximum
157 kernel density and 218.99 km using the geographic centroid of window predictions (Figure 6;
158 Table S2). Mean predicted locations across all windows consistently separate populations in the
159 Americas, West Africa, East Africa, southeast Asia, and Papua New Guinea; consistent with the
160 major population subdivisions described via PCA in Pearson et al. (2019). We also see good
161 discrimination within clusters, particularly in southeast Asia where the average test error is less
162 than 200km for all but two localities. Error is highest in West Africa, where mean predictions tend
163 towards the center of a set of regional collecting localities (Figure 6). These patterns are consistent
164 with previous findings of fine-scale spatial structure in *P. falciparum* in Cambodia (Miotto et al.,
165 2013) and low levels of relative genetic differentiation (as measured by F_{ST}) in Africa (Pearson
166 et al., 2019).

167 Rates of mixed-strain infection are elevated in West Africa relative to Southeast Asia (Zhu
168 et al., 2019; Pearson et al., 2019), which we hypothesized could explain the higher prediction error
169 in this region. To test this effect we plotted *Locator*'s centroid prediction error as a function of
170 within-host diversity (F_{WS} ; Auburn et al. (2012)). F_{WS} measures the proportion of population
171 genetic diversity present in individual hosts, with a value of 0 representing maximum within-
172 host diversity and 1 minimum within-host diversity. If mixed-strain infections explain outliers of
173 prediction error, we would expect that samples with the highest prediction error had low F_{WS} .
174 Instead we found a weak positive relationship (Figure S6), with the highest prediction errors seen
175 in samples with maximum F_{WS} (i.e., minimum infection diversity). Test error then likely reflects
176 low levels of differentiation within *Plasmodium* lineages in West Africa rather than local prevalence
177 of mixed-strain infections.

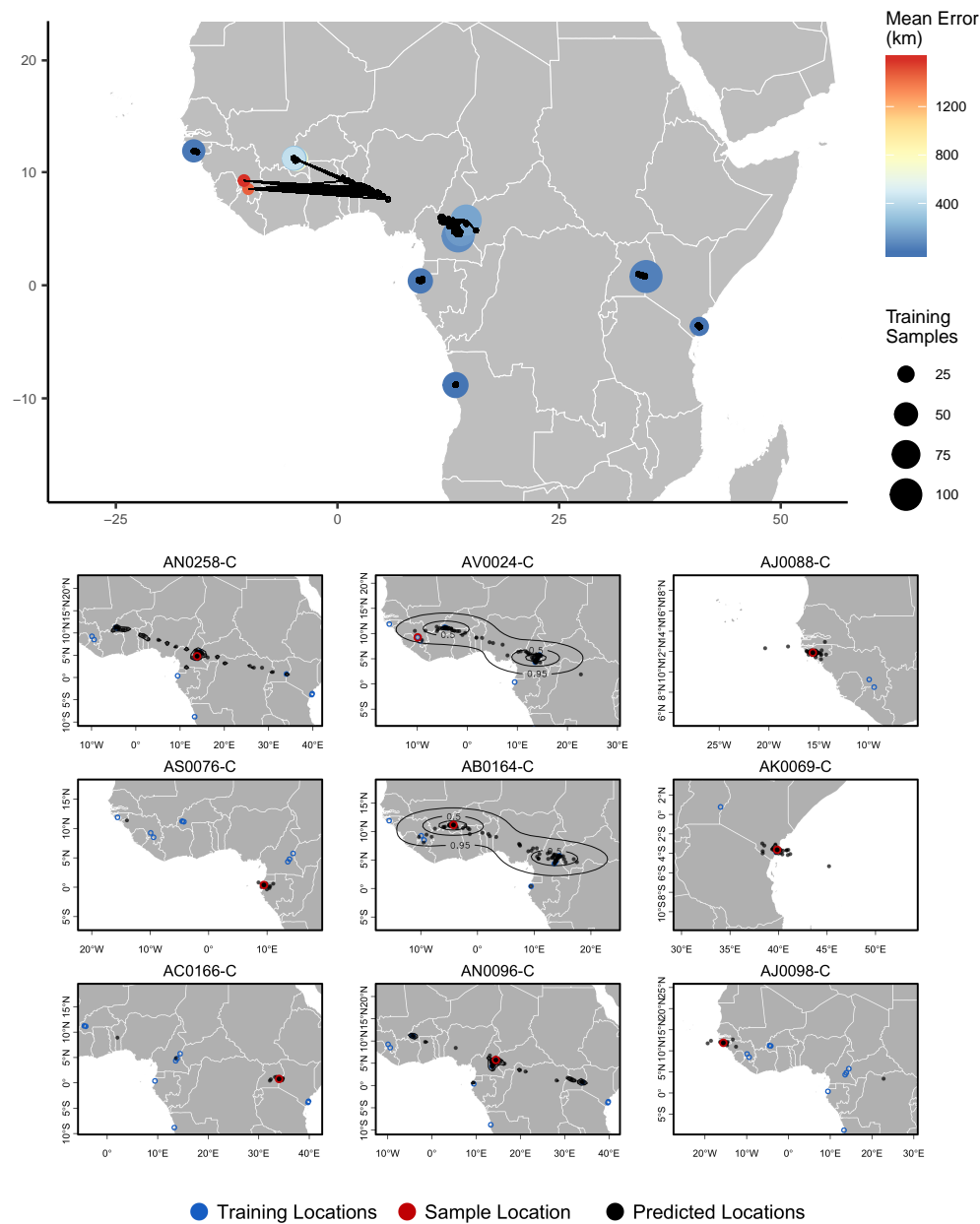


Figure 5: **Top** – Predicted locations for 153 *Anopheles gambiae* / *coluzzii* genomes from the AG1000G panel, using 612 training samples and a 2Mbp window size. The geographic centroid of per-window predictions for each individual is shown in black points, and lines connect predicted to true locations. Sample localities are colored by the mean test error with size scaled to the number of training samples. **Bottom** – Uncertainty from predictions in 2Mbp windows. Contours show the 95%, 50%, and 10% quantiles of a two-dimensional kernel density across windows.

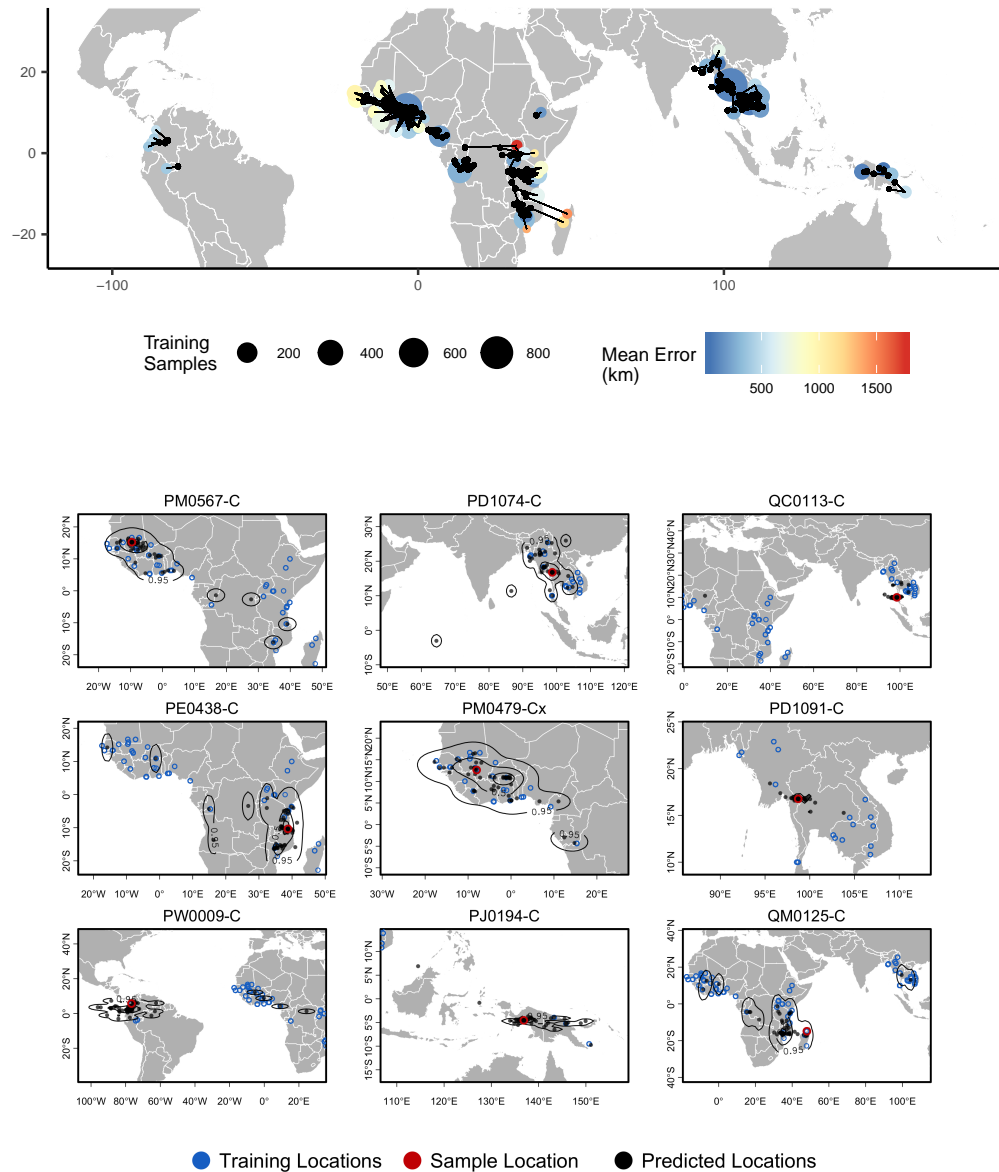


Figure 6: **Top** – Predicted locations for 881 *Plasmodium falciparum* from the *Plasmodium falciparum* Community Project (Pearson et al., 2019) (5% of samples for each collecting locality), using 5084 training samples and a 500Kbp window size. The geographic centroid of per-window predictions for each individual is shown in black points, and lines connect predicted to true locations. Sample localities are colored by the mean test error with size scaled to the number of training samples. **Bottom** – Uncertainty from predictions in 500Kbp windows. Contours show the 95%, 50%, and 10% quantiles of a two-dimensional kernel density across windows.

178 Again we found that visualizing per-window predictions reflects expected patterns of uncer-
179 tainty in samples with high mean prediction error. For example, sample QM0215-C was collected
180 in Madagascar and has a mean predicted location in Mozambique, but the spread of per-window
181 predictions indicates a 95% interval that includes the true locality (Figure 6, bottom right).

182 The good performance we observed on this dataset also highlights a strength of *Locator*'s
183 model-free approach. Recall that the sequencing strategy of preparing libraries from human blood
184 samples suggests variant calls represent binned allele frequencies across the population of *Plas-*
185 *modium* in a human blood sample rather than SNPs in a single *Plasmodium* individual. From the
186 perspective of the network; however, the input genotypes are simply a set of normalized vectors,
187 and the network can approximate the relationship between these vectors and the spatial location
188 of training samples regardless of the generative process.

189 Human Populations

190 For humans in the HGDP dataset, the location with highest kernel density across all windows has
191 a median test error of 85km, and the centroid of window predictions has a median error of 452.6
192 kilometers (Figure 7, Table S2). Visualizing the geographic distribution of predictions across the
193 genome shows that predictions tend to cluster around the true reported sampling location, but
194 also extend towards other sampling locations in a manner that reflects known patterns of human
195 migration.

196 For example, the two largest individual errors in our analysis are found in South African
197 Bantu individuals and Xibo people from western China. Predicted locations of South African
198 Bantu people project towards the historic source of Bantu migrations in west Africa (De Filippo
199 et al., 2012), with some regions of the genome also projecting in the direction of east African
200 Bantu populations (Figure 7, sample HGDP00993). In the case of Xibo people from western
201 China *Locator* consistently predicts locations in Manchuria, central China, and southern Siberia
202 – significantly east of the true sample location. This may reflect the known movement of this
203 population, which historically originated in Manchuria and was resettled in western China during
204 the 18th century (Gorelova, 2002; Zikmundová, 2013) (Figure 7, sample HGDP01250). A sample
205 of individual-level predictions is included in Figure 7.

206 To test whether outlier geographic predictions reflect error in the model fitting procedure
207 versus true variation in ancestry in a given region of the genome, we ran principal component
208 analyses on windows for which a Maya individual (sample HGDP00871) has predicted locations
209 in Europe and Africa. In these windows the Maya sample clusters with other individuals from the
210 regions predicted by *Locator* – western Europe and Africa, respectively – rather than with other
211 individuals from the Americas (Figure S9). This suggests outlier predictions reflect variation in
212 ancestry in different regions of the genome, rather than stochastic error in model fitting.

213 We also examined how recombination rate interacts with the accuracy of *Locator* predictions
214 generated from different regions of the genome. We might expect recombination rate to affect
215 accuracy because in regions of the genome with higher recombination, there are a greater number
216 of distinct genealogies, and hence a given sample has inherited from a larger subset of the possible
217 ancestors. Test error was estimated as the distance in kilometers from the true sampling loca-
218 tion to the geographic centroid of the cloud of per-window predictions, and is shown in figure 8
219 plotted against local recombination rates from the HapMap genetic map (International HapMap
220 Consortium, 2003). We find a relatively strong negative correlation ($p < 0.0001$, $R^2 = 0.27$) –
221 windows with the lowest recombination rates in general have the highest prediction error. This
222 is consistent with our expectation that regions of the genome representing a greater number of
223 marginal genealogies will yield more accurate predictions of a sample's location.

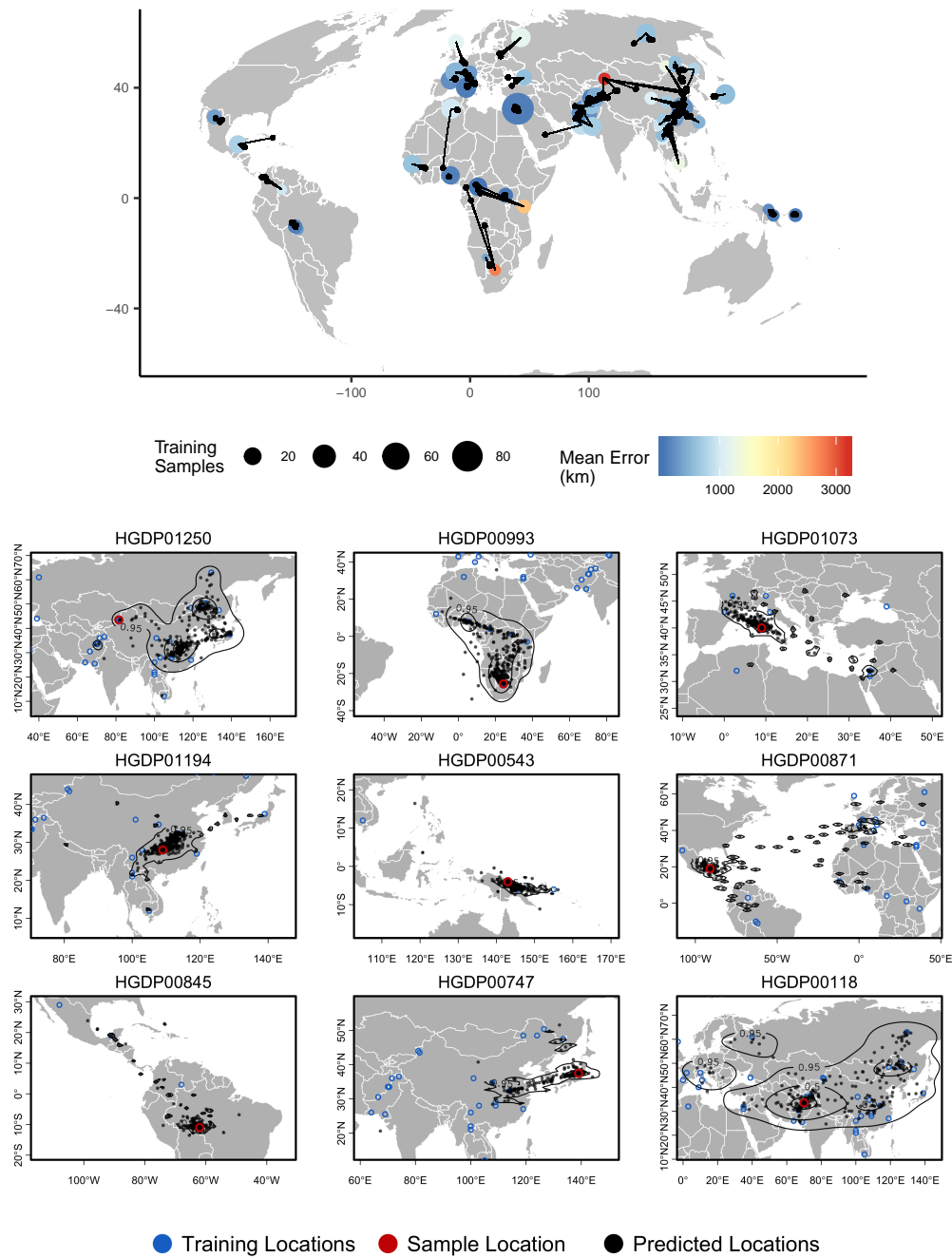


Figure 7: Top – Predicted locations for 162 individuals from the HGDP panel, using 773 training samples and a 10Mbp window size. The geographic centroid of per-window predictions for each individual is shown in black points, and lines connect predicted to true locations. Sample localities are colored by the mean test error with size scaled to the number of training samples. Bottom – Uncertainty from predictions in 10Mbp windows. Contours show the 95%, 50%, and 10% quantiles of a two-dimensional kernel density across windows.

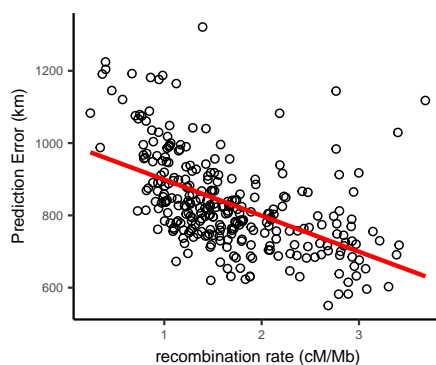


Figure 8: Per-window test error and mean recombination rate for human populations in the HGDP dataset. The top 2% of windows by test error were excluded from this analysis. The slope of the least-squares linear fit is -99.9723 km/(cM/Mbp) and has adjusted $R^2 = 0.2704$.

224 Effects of Unsampled Populations

225 In figures S7 and S8 we show predictions from a single window in the *Anopheles* dataset when
226 no samples from a given geographic region are included in the training set at two scales – either
227 dropping only sites from a specific sampling location or dropping all sites from a given country.
228 Prediction error is much higher for individuals from regions excluded from training – increasing
229 from a median of 14km when training and test samples are randomly split to 116km when excluding
230 individual localities, and 778km when excluding whole countries.

231 In most cases predicted locations appear to project towards the nearest locality included in the
232 training set (figure S8). This is particularly the case when populations at the edge of the map
233 are excluded. **Locator** networks appear to learn something about the boundaries of the landscape
234 based on the distribution of training points, and show a tendency to project towards the middle of
235 the landscape when given a small number of SNPs (e.g., the top right panel of Figure 2A), a trivial
236 optimization of the loss function. We also see evidence of **Locator** learning some nonlinear aspects
237 of population structure in the sample. For example, when Angolan *A. coluzzii* are excluded from
238 the training set many of their predicted locations project towards the *A. coluzzii* sample localities
239 in Burkina Faso rather than the much closer sampling localities for *A. gambiae* in Cameroon and
240 Gabon. In general we find that **Locator** can interpolate unsampled localities relatively well when
241 genetic differentiation is smooth over the landscape (as among *A. gambiae* localities in west Africa),
242 but does not extrapolate outside the bounds of the training set. Sampling the full landscape, or
243 at least a sufficient portion thereof, is thus an important consideration in running our method.

244 Discussion

245 The correlation of genealogy and geography leaves genetic signals of ancestral location across the
246 genome that one can leverage for practical inference. For instance, tracking the migratory routes
247 of disease vectors such as *Anopheles* (Huestis et al., 2019) could in principle be achieved if one
248 could accurately predict origin from DNA sequence data. Similarly, establishing the location of
249 origin from biological samples is critical to anti-poaching conservation efforts (Wasser et al., 2004).
250 In this report we present a new tool, **Locator**, which uses a deep neural network to predict the
251 geographic location of a sample on the basis of its genotype. We show that **Locator** is highly
252 accurate, computationally efficient, and can scale to thousands of genomes.

253 In simulations we showed that our method returns the same results as a state-of-the-art model-

254 based approach, SPASIBA (Guillot et al., 2015), and does so at least an order of magnitude faster.
255 We show that the accuracy of our estimator is naturally measured in terms of the dispersal rate of
256 the population and that predictions from **Locator** are consistently within 3–4 generations of mean
257 dispersal across a wide range of dispersal distances (Figure 2, Figure S2). However we found that
258 the greatest increase in performance relative to the model-based approach is in empirical data for
259 which the assumption of smooth variation in allele frequencies across the landscape is unlikely to
260 hold, such as the complex multi-species *Anopheles* sample analyzed here (Figure 3).

261 **Locator**'s computational efficiency makes it practical to estimate uncertainty through resam-
262 pling approaches like windowed analysis or bootstrapping over the complete genotype matrix. The
263 full windowed analysis of the HGDP data took roughly 30 hours to run on a single GPU, and
264 windowed analysis of all 5,965 complete *Plasmodium* genomes took just 8 hours. Thus training
265 **Locator** models for biobank-scale datasets including whole genomes of tens or hundreds of thou-
266 sands of samples is well within reach, particularly if windows can be run on separate GPUs. This
267 allows us to estimate uncertainty in predicted locations due both to our prediction methodology as
268 well as biology; with repeated training runs integrating over error associated with network training
269 and prediction and the windowed analysis allowing us to predict geographic origins for regions of
270 the genome reflecting distinct sets of genealogical relationships.

271 Disentangling these sources of error is challenging, but analysis of human data for which we
272 have strong prior knowledge of recent population movements suggests that much of the variation in
273 genome-wide prediction we see reflects historic patterns of migration rather than simple prediction
274 error. For example, genomes from Hazara individuals in central Asia return predicted locations
275 extending from central Asia to Mongolia (Figure 7 bottom, sample HGDP00118), which is consis-
276 tent with historic records (Qamar et al., 2002), previous analysis of Y chromosome data (Zerjal
277 et al., 2003), and identity-by-descent tract sharing (Lawson et al., 2012) all of which find evidence
278 of recent shared ancestry between Mongolian and Hazara individuals. Similarly some Maya indi-
279 viduals found to have a small proportion of European ancestry in previous analyses (Rosenberg
280 et al., 2002) have predicted locations extending from central Mexico across the Atlantic to Europe
281 and west Africa in windowed **Locator** analysis (Figure 7 bottom, sample HGDP00871), and these
282 signals are replicated in principal components analysis (Figure S9).

283 This also points to a critical consideration in running any form of supervised population clus-
284 tering. Information about population structure comes only from the relative relationships among
285 training and test samples, and interpretations can only be made relative to the set of training
286 samples used. In the case of the HGDP panel, samples were intentionally selected to cover what
287 were thought to be distinctive populations reflecting a vaguely pre-modern distribution of human
288 genetic diversity (Harry and Marks, 1999), and so would probably not be a good reference set
289 for random individuals drawn from regions or groups with recent histories of large population
290 movements such as the United States.

291 Here we have shown that our method, **Locator**, is fast, accurate, and scales well to large
292 samples. However we see several next steps that could improve the approach. First, our current
293 implementation uses only diploid genotypes and does not pass the network any direct information
294 about haplotype structure (though in theory the fully-connected nature of our network could allow
295 inference of pairwise correlations among sites). Incorporating SNP position information and phased
296 haploid sequences would likely increase inferential power, as in the case of unsupervised clustering
297 (Lawson et al., 2012). Second, our network currently uses a simple fully-connected architecture; it
298 could be that other network architectures such as recurrent neural networks might be better suited
299 for this task (e.g., Adrion et al., 2019). Indeed the application of deep learning to population
300 genetics is still in its infancy and we imagine much progress will be made in the coming years
301 along these lines.

302 Methods

303 Preprocessing

304 **Locator** transforms input data in VCF or Zarr format to vectors of derived allele counts per
305 individual using the scikit-allel (Miles and Harding, 2017) and numpy (Van Der Walt et al., 2011)
306 libraries. Sites with missing data are replaced with two draws from a binomial distribution with
307 probability equal to the frequency of the derived allele across all individuals – a discrete version of
308 the common practice of assigning missing data as the mean allele frequency in genotype PCAs (e.g.
309 the default settings for PCA in the R package adegenet (Jombart, 2008)). We provide functions
310 for filtering SNPs based on minor allele count, and by default remove singleton sites from the
311 alignment prior to model fitting. The geographical x and y coordinates are scaled to have mean
312 0 and variance 1 prior to training, while allele counts are scaled prior to model fitting by a batch
313 normalization layer within the network. Batch normalization Z-normalizes activations of a neural
314 network during training to reduce shifts in the distribution of parameter values across batches,
315 which allows faster learning rates and sometimes reduces overfitting (Ioffe and Szegedy, 2015).

316 **Locator** selects a user-defined fraction of the samples with known locations to use in training the
317 model (the default is 0.9); remaining samples with known locations are kept aside as “validation”
318 samples. The validation set is used to tune the learning rate of the optimizer and set the stopping
319 time of model training, but does not directly contribute to the loss used to fit model parameters.
320 Throughout this manuscript we use “validation loss” to refer to error estimated on the validation
321 set, and “test error” to refer to error calculated on a set of samples entirely held out from the
322 model training procedure.

323 Network

324 We use the unphased, diploid genotype vector of each individual as input to the network, whose
325 target output is the two-dimensional coordinates of that individual in space. **Locator** uses a deep
326 neural network consisting of a stack of fully-connected “dense” layers, implemented using the Keras
327 (Chollet et al., 2015) frontend to tensorflow (Abadi et al., 2015). Roughly speaking, the network
328 is trained to estimate a nonlinear function mapping genotypes to locations using gradient-based
329 optimization. Models start with randomized initial parameters and are fit to data by looping
330 through the training set and iteratively adjusting the weights and biases of the network. We use
331 an early stopping function to monitor loss during training and under default settings stop training
332 runs when validation loss has not improved for 100 epochs. We also use a learning rate scheduler
333 to decrease the learning rate of the optimizer when validation loss stops improving, which we found
334 to be effective in preventing the trajectories of training and validation loss from diverging. The
335 program also outputs a plot of training and validation loss after each training run (Figure S1).

336 **Locator**’s architecture uses a batch normalization layer followed by a sequence of fully-connected
337 layers with a dropout layer in the middle of the network (Figure 1). The “dropout” layer sets a
338 random selection of weights to zero during each training step, which helps prevent overfitting (Sri-
339 vastava et al., 2014). Our implementation allows users to adjust the shape of the network, but
340 current default settings use 10 dense layers of 256 nodes each with “ELU” activations (Clevert
341 et al., 2015) and a 25% dropout after the fifth layer. We describe performance under varying
342 network width and depth in Supplementary Figure S3. In general we found that all networks with
343 over four layers perform similarly.

We use the Adam optimizer (Kingma and Ba, 2014) with Euclidean distance as a loss function:

$$\text{loss} = \sqrt{(x_{\text{predicted}} - x_{\text{true}})^2 + (y_{\text{predicted}} - y_{\text{true}})^2}. \quad (1)$$

344 Uncertainty and Genome-wide Variation

345 Individuals are born at a single location, but have inherited their genomes as a mosaic from
346 ancestors spreading geographically into the past (as discussed in, for instance, Wright (1943);
347 Kelleher et al. (2016); Bradburd and Ralph (2019)). Any signal our method hopes to extract from
348 the data must be due to geographic signal of recent ancestors shared between the test and training
349 datasets. This suggests that any analogous method must quantify, roughly, “which modern day
350 populations are most similar to this genome?”. The spatial spread of genetic relatedness both
351 back in time from an individual’s to its ancestors’ locations and forward in time from ancestors to
352 the present-day location of training samples means that even a perfect inference algorithm should
353 have significant uncertainty associated with any predicted location from genetic data, and the
354 magnitude of uncertainty should be in part a function of the dispersal rate of the population. In
355 particular, no such method can infer locations more accurately than the mean dispersal distance,
356 because in most cases an individual’s genome is not informative about where they live relative to
357 their parents. Besides this fundamental limit to uncertainty, error in georeferencing of training
358 samples and in model fitting will introduce additional prediction uncertainty.

359 We use a windowed analysis across the genome to describe this uncertainty, which is possible
360 thanks to *Locator*’s computational efficiency. Genealogical relatedness on each contiguous stretch
361 of genome can be described by a sequence of genealogical trees, separated by ancestral recombina-
362 tion events. By running *Locator* on a particular window of the genome, we restrict inference to
363 a subset of these marginal trees, and hence to a subset of the genetic relationships between test
364 and training samples. Predictions from different regions of the genome can then be visualized as
365 a cloud of points, and the distribution of these points in space gives us a rough idea of the uncer-
366 tainty associated with an individual-level prediction. Because windowed analyses involve repeated
367 training runs from randomized starting parameters, they also help us to integrate over uncertainty
368 associated with the model fitting process.

369 Some datasets lack the size or reference alignments necessary to conduct windowed analyses.
370 In this case we recommend uncertainty be assessed by training replicate models on bootstrapped
371 samples drawn from a single set of unlinked SNPs (that is, resampling SNPs with replacement).
372 Though this procedure does not reduce the number of marginal trees represented in the data, it
373 does allow us to assess uncertainty associated with model training and prediction. In both cases
374 we summarize uncertainty in predicted locations by estimating a two-dimensional kernel density
375 surface over a set of predicted locations, and provide plotting scripts to visualize the 95%, 50%,
376 and 10% quantiles in geographic space (see figures 5–7 for examples). The location of an individual
377 can then be predicted as either the location with highest kernel density (the modal prediction) or
378 the geographic center of the cloud of predictions (the mean prediction).

379 We tested this approach in simulated data and in all empirical datasets. To explore factors
380 affecting the accuracy of predicted locations generated from different regions of the genome, we
381 also examined the relationship between recombination rate and test error from windowed *Locator*
382 runs on human data from the HGDP panel (Bergström et al., 2019). Recombination rates for
383 each window were estimated by averaging per-base rates from the HapMap project (International
384 HapMap Consortium, 2003).

385 Simulations

386 We first evaluated our method on genotypes from populations simulated by SLiM v3 (Haller and
387 Messer, 2019), using the model of continuous space described in Batthey et al. (2019). We simulated
388 a 50×50 unit square landscape with expected density (d) of 5 individuals per unit area, resulting
389 in census sizes of around 12,500. We varied the mean parent-offspring dispersal distance σ across
390 simulations from 0.45 to 3, to create populations with varying levels of isolation by distance.
391 In terms of Wright’s “neighborhood size” (Wright, 1946), defined as $N_{loc} = 4\pi\sigma^2d$, this yields

392 populations with neighborhood sizes from 13 to 565. Each diploid individual carried two copies of
393 a 10^8 bp chromosome on which mutations and recombinations occurred at a rate of 10^{-8} per bp
394 per generation. Simulations were run until all extant individuals shared a single common ancestor
395 within the simulation at all locations on the genome (i.e., the tree sequence had coalesced). 500
396 individuals were randomly sampled from the final generation of each simulation for use in model
397 fitting.

398 We selected 50 individuals from each simulation as a validation set and ran `Locator` while
399 varying the number of training samples from 10 to 450 and the number of SNPs from 100 to 100,000.
400 The SNPs used were a subset sampled from the full genotype matrix without replacement and thus
401 mimic the semi-random distribution of genome-wide SNPs generated by reduced-representation
402 sequencing approaches like RADseq (Etter et al., 2012). To compare performance with an existing
403 model-based approach, we also ran SPASIBA (Guillot et al., 2015) on the simulation with $\sigma = 0.44$
404 using 450 training samples and varying the number of SNPs from 100 to 100,000. `Locator` was
405 run on a CUDA-enabled GPU and SPASIBA was run on 80 CPU cores. Last, we ran a windowed
406 analysis on the $\sigma = 0.63$ (neighborhood size ≈ 25) simulation in `Locator` using a 2Mbp window
407 size (each window then contains $\approx 8,000$ SNPs).

408 Empirical Data

409 We applied `Locator` to three whole-genome resequencing datasets of geographically widespread
410 samples: (1) 765 mosquitoes from the *Anopheles gambiae* / *coluzzii* species complex collected
411 across sub-Saharan Africa (Miles et al., 2017), (2) 5,965 samples of the malaria parasite *Plasmodium*
412 *falciparum* sequenced from human blood samples collected across Papua New Guinea, southeast
413 Asia, sub-Saharan Africa, and northern South America (Pearson et al., 2019) and (3) whole-genome
414 data for 56 human populations from the Human Genome Diversity Project (Bergström et al., 2019).
415 Genotype calls for the *Anopheles* dataset are available at <https://www.malariagen.net/data/ag1000g-phase1-ar3>, for *P. falciparum* at <https://www.malariagen.net/resource/26>, and for
416 human data at <ftp://ngs.sanger.ac.uk/production/hgdp>. We used VCF files as provided with
417 no further postprocessing.
418

419 The *Plasmodium falciparum* dataset is unusual relative to our other empirical examples in
420 that sequencing libraries were prepared from blood samples without filtering for coinfections or
421 isolating individual *Plasmodium*. Sequence reads returned from short read sequencing then reflect
422 the population of *Plasmodium* present in a human blood sample, or even multiple lineages of
423 parasite if an individual is co-infected with multiple strains (Zhu et al., 2019), rather than individual
424 *Plasmodium*. The VCFs we analyzed were prepared by aligning illumina short read sequences to
425 the *Plasmodium falciparum* reference genome prepared by the Pf3K project (Pf3K Consortium
426 (2016); <https://www.malariagen.net/data/pf3K-5>), then calling SNPs in GATK (McKenna
427 et al., 2010). Variant calls then represent the pool of mutations present in the infecting *Plasmodium*
428 population rather than SNPs in a single individual. We used only field-collected samples from the
429 “analysis” set, as described in (Pearson et al., 2019).

430 For the *Anopheles* dataset we ran `Locator` in 2Mbp windows across the genome with a randomly
431 selected 10% of individuals held out as a test set. We also ran SPASIBA on subsets sampled from
432 the first five million base pairs of chromosome 2L while varying the number of SNPs from 100 to
433 100,000. For the *P. falciparum* dataset we used 500kb windows and held out 5% of samples from
434 each collection locality as a test set. Last, for humans we used 10Mbp windows and selected three
435 individuals from each HGDP population to hold out as a test set. Window sizes in each case were
436 chosen to include roughly 100,000–200,000 SNPs per window. All empirical analyses were run
437 with default settings (10×256 network size, patience 100, 25% dropout, a random 10% of training
438 samples used for validation).

439 We also tested `Locator`’s performance with empirical data when the true location is not rep-

resented in the training sample. To do this we ran a series of models on 10,000 SNPs randomly selected from the first 5Mbp of chromosome 2L in the *Anopheles* data. For each run we held out all samples from a given sampling locality from the training set, then predicted the locations of these individuals using the trained model. We also tested this approach while holding out all samples collected in a given country, which eliminates even nearby localities from the training set.

Data & Code

Locator is implemented as a command-line program written in Python: www.github.com/kern-lab/locator. SNP calls for the *Anopheles* dataset are available at <https://www.malariagen.net/data/ag1000g-phase1-ar3>, for *P. falciparum* at <https://www.malariagen.net/resource/26>, and for the HGDP at <ftp://ngs.sanger.ac.uk/production/hgdp>. This publication uses data from the MalariaGEN Plasmodium falciparum Community Project as described in Pearson et al. (2019). Statistical analyses and many plots were produced in R (R Core Team, 2018).

Competing Interests

The authors declare no competing interests.

Acknowledgements

We thank members of the Kern-Ralph co-lab, Daniel Schrider, Matthew Hahn, and Ethan Linck for comments and suggestions on this work, and Mara Lawniczak for the suggestion to look at the *Plasmodium* dataset. CJB and ADK were funded by NIH award R01GM117241. PLR was funded in part by an I3 award from the University of Oregon.

References

- 459
- 460 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.
461 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, An-
462 drew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser,
463 Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray,
464 Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul
465 Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden,
466 Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale
467 machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>. Software
468 available from tensorflow.org.
- 469 Jeffrey R Adrion, Jared G Galloway, and Andrew D Kern. Inferring the landscape of recombination
470 using recurrent neural networks. *bioRxiv*, page 662247, 2019.
- 471 Sarah Auburn, Susana Campino, Olivo Miotto, Abdoulaye A Djimde, Issaka Zongo, Magnus
472 Manske, Gareth Maslen, Valentina Mangano, Daniel Alcock, Bronwyn MacInnis, et al. Charac-
473 terization of within-host *Plasmodium falciparum* diversity using next-generation sequence data.
474 *PloS one*, 7(2):e32891, 2012.
- 475 Yael Baran, Inés Quintela, Ángel Carracedo, Bogdan Pasaniuc, and Eran Halperin. Enhanced
476 localization of genetic samples through linkage-disequilibrium correction. *The American Journal*
477 *of Human Genetics*, 92(6):882–894, 2013.
- 478 CJ Batthey, Ethan B Linck, Kevin L Epperly, Cooper French, David L Slager, Paul W Sykes Jr,
479 and John Klicka. A migratory divide in the Painted Bunting (*Passerina ciris*). *The American*
480 *Naturalist*, 191(2):259–268, 2018.
- 481 CJ Batthey, Peter L Ralph, and Andrew D Kern. Space is the place: Effects of continuous spatial
482 structure on analysis of population genetic data. *BioRxiv*, page 659235, 2019.
- 483 Anders Bergström, Shane A. McCarthy, Ruoyun Hui, Mohamed A. Almarri, Qasim Ayub, Petr
484 Danecek, Yuan Chen, Sabine Felkel, Pille Hallast, Jack Kamm, Hélène Blanché, Jean-François
485 Deleuze, Howard Cann, Swapan Mallick, David Reich, Manjinder S. Sandhu, Pontus Skoglund,
486 Aylwyn Scally, Yali Xue, Richard Durbin, and Chris Tyler-Smith. Insights into human genetic
487 variation and population history from 929 diverse genomes. *bioRxiv*, 2019. doi: 10.1101/674986.
488 URL <https://www.biorxiv.org/content/early/2019/06/27/674986>.
- 489 Anand Bhaskar, Adel Javanmard, Thomas A Courtade, and David Tse. Novel probabilistic models
490 of spatial genetic ancestry with applications to stratification correction in genome-wide associ-
491 ation studies. *Bioinformatics*, 33(6):879–885, 2016.
- 492 Gideon S. Bradburd and Peter L. Ralph. Spatial population genetics: It’s about
493 time. *Annual Review of Ecology, Evolution, and Systematics*, 50(1):427–449, 2019.
494 doi: 10.1146/annurev-ecolsys-110316-022659. URL <https://doi.org/10.1146/annurev-ecolsys-110316-022659>.
- 495
- 496 Natalie Breidenbach, Oliver Gailing, and Konstantin V Krutovsky. Assignment of frost tolerant
497 coast redwood trees of unknown origin to populations within their natural range using nuclear
498 and chloroplast microsatellite genetic markers. *bioRxiv*, page 732834, 2019.
- 499 Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.

- 500 Jeffrey Chan, Valerio Perrone, Jeffrey Spence, Paul Jenkins, Sara Mathieson, and Yun Song.
501 A likelihood-free inference framework for population genetic data using exchangeable neural
502 networks. In *Advances in Neural Information Processing Systems*, pages 8594–8605, 2018.
- 503 François Chollet et al. Keras. <https://keras.io>, 2015.
- 504 Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network
505 learning by exponential linear units (ELUs). *arXiv preprint arXiv:1511.07289*, 2015.
- 506 Qian Cong, Jinhui Shen, Jing Zhang, Wenlin Li, Lisa N. Kinch, John V. Calhoun, Andrew D.
507 Warren, and Nick V. Grishin. Genomics reveals the origins of ancient specimens. *bioRxiv*, 2019.
508 doi: 10.1101/752121. URL <https://www.biorxiv.org/content/early/2019/09/04/752121>.
- 509 Cesare De Filippo, Koen Bostoen, Mark Stoneking, and Brigitte Pakendorf. Bringing together
510 linguistic and genetic evidence to test the Bantu expansion. *Proceedings of the Royal Society B:
511 Biological Sciences*, 279(1741):3256–3263, 2012.
- 512 Arun Durvasula and Sriram Sankararaman. A statistical model for reference-free inference of
513 archaic local ancestry. *PLoS genetics*, 15(5):e1008175, 2019.
- 514 Paul D Etter, Susan Bassham, Paul A Hohenlohe, Eric A Johnson, and William A Cresko. SNP
515 discovery and genotyping for evolutionary genetics using RAD sequencing. In *Molecular methods
516 for evolutionary genetics*, pages 157–178. Springer, 2012.
- 517 Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The unreasonable effectiveness of convolu-
518 tional neural networks in population genetic inference. *Molecular biology and evolution*, 36(2):
519 220–238, 2018.
- 520 Liliya M Gorelova. *Manchu grammar*. Brill, 2002.
- 521 Gilles Guillot, Hákon Jónsson, Antoine Hinge, Nabil Manchih, and Ludovic Orlando. Accurate
522 continuous geographic assignment from low- to high-density SNP data. *Bioinformatics*, 32
523 (7):1106–1108, 11 2015. ISSN 1367-4803. doi: 10.1093/bioinformatics/btv703. URL <https://doi.org/10.1093/bioinformatics/btv703>.
- 525 Benjamin C Haller and Philipp W Messer. SLiM 3: Forward genetic simulations beyond the
526 Wright–Fisher model. *Molecular biology and evolution*, 36(3):632–637, 2019.
- 527 Debra Harry and Jonathan Marks. Human population genetics versus the HGDP. *Politics and
528 the Life Sciences*, 18(2):303–305, 1999.
- 529 Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document
530 analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- 531 Diana L. Huestis, Adama Dao, Moussa Diallo, Zana L. Sanogo, Djibril Samake, Alpha S. Yaro,
532 Yossi Ousman, Yvonne-Marie Linton, Asha Krishna, Laura Veru, Benjamin J. Krajacich, Roy
533 Faiman, Jenna Florio, Jason W. Chapman, Don R. Reynolds, David Weetman, Reed Mitchell,
534 Martin J. Donnelly, Elijah Talamas, Lourdes Chamorro, Ehud Strobach, and Tovi Lehmann.
535 Windborne long-distance migration of malaria mosquitoes in the Sahel. *Nature*, 2019. doi:
536 10.1038/s41586-019-1622-4. URL <https://doi.org/10.1038/s41586-019-1622-4>.
- 537 International HapMap Consortium. The international HapMap project. *Nature*, 426(6968):789,
538 2003.
- 539 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
540 reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

- 541 Thibaut Jombart. adegenet: a R package for the multivariate analysis of genetic markers. *Bioin-*
542 *formatics*, 24(11):1403–1405, 2008.
- 543 Jerome Kelleher, AM Etheridge, Amandine Véber, and Nicholas H Barton. Spread of pedigree
544 versus genetic ancestry in spatially distributed populations. *Theoretical population biology*, 108:
545 1–12, 2016.
- 546 Andrew D Kern and Daniel R Schrider. diploS/HIC: an updated approach to classifying selective
547 sweeps. *G3: Genes, Genomes, Genetics*, 8(6):1959–1970, 2018.
- 548 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
549 *arXiv:1412.6980*, 2014.
- 550 Daniel John Lawson, Garrett Hellenthal, Simon Myers, and Daniel Falush. Inference of population
551 structure using dense haplotype data. *PLoS genetics*, 8(1):e1002453, 2012.
- 552 Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew
553 Kernysky, Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, et al. The genome
554 analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.
555 *Genome research*, 20(9):1297–1303, 2010.
- 556 Alistair Miles and Nick Harding. cggh/scikit-allel: v1.1.8, July 2017. URL [https://doi.org/10.](https://doi.org/10.5281/zenodo.822784)
557 [5281/zenodo.822784](https://doi.org/10.5281/zenodo.822784).
- 558 Alistair Miles, Nicholas J Harding, and the AG1000G consortium. Genetic diversity of the African
559 malaria vector *Anopheles gambiae*. *Nature*, 552(7683):96, 2017.
- 560 Olivo Miotto, Jacob Almagro-Garcia, Magnus Manske, Bronwyn MacInnis, Susana Campino,
561 Kirk A Rockett, Chanaki Amaratunga, Pharath Lim, Seila Suon, Sokunthea Sreng, et al. Multi-
562 ple populations of artemisinin-resistant *Plasmodium falciparum* in Cambodia. *Nature genetics*,
563 45(6):648, 2013.
- 564 Mehreen R Mughal and Michael DeGiorgio. Localizing and classifying adaptive targets with trend
565 filtered regression. *Molecular biology and evolution*, 36(2):252–270, 2018.
- 566 Richard D Pearson, Roberto Amato, Dominic P Kwiatkowski, and MalariaGEN Plasmodium fal-
567 ciparum Community Project. An open dataset of *Plasmodium falciparum* genome variation in
568 7,000 worldwide samples. *bioRxiv*, 2019. doi: 10.1101/824730. URL [https://www.biorxiv.](https://www.biorxiv.org/content/early/2019/11/07/824730)
569 [org/content/early/2019/11/07/824730](https://www.biorxiv.org/content/early/2019/11/07/824730).
- 570 Pf3K Consortium. The Pf3K project (2016): pilot data release 5, 2016. URL [www.malariagen.](http://www.malariagen.net/data/pf3k-5)
571 [net/data/pf3k-5](http://www.malariagen.net/data/pf3k-5).
- 572 Pierre Pudlo, Jean-Michel Marin, Arnaud Estoup, Jean-Marie Cornuet, Mathieu Gautier, and
573 Christian P Robert. Reliable ABC model choice via random forests. *Bioinformatics*, 32(6):
574 859–866, 2015.
- 575 Raheel Qamar, Qasim Ayub, Aisha Mohyuddin, Agnar Helgason, Kehkashan Mazhar, Atika Man-
576 soor, Tatiana Zerjal, Chris Tyler-Smith, and S Qasim Mehdi. Y-chromosomal DNA variation in
577 Pakistan. *The American Journal of Human Genetics*, 70(5):1107–1124, 2002.
- 578 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for
579 Statistical Computing, Vienna, Austria, 2018. URL <https://www.R-project.org/>.
- 580 John Michael Rañola, John Novembre, and Kenneth Lange. Fast spatial ancestry via flexible allele
581 frequency surfaces. *Bioinformatics*, 30(20):2915–2922, 2014.

- 582 Noah A Rosenberg, Jonathan K Pritchard, James L Weber, Howard M Cann, Kenneth K Kidd,
583 Lev A Zhitovskiy, and Marcus W Feldman. Genetic structure of human populations. *science*,
584 298(5602):2381–2385, 2002.
- 585 Daniel R Schrider and Andrew D Kern. S/HIC: Robust identification of soft and hard sweeps using
586 machine learning. *PLoS Genet*, 12(3):e1005928, Mar 2016. doi: 10.1371/journal.pgen.1005928.
- 587 Daniel R Schrider, Julien Ayroles, Daniel R Matute, and Andrew D Kern. Supervised machine
588 learning reveals introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLoS*
589 *genetics*, 14(4):e1007341, 2018.
- 590 Sara Sheehan and Yun S Song. Deep learning for population genetic inference. *PLoS computational*
591 *biology*, 12(3):e1004845, 2016.
- 592 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
593 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*
594 *learning research*, 15(1):1929–1958, 2014.
- 595 Lauren Alpert Sugden, Elizabeth G Atkinson, Annie P Fischer, Stephen Rong, Brenna M Henn,
596 and Sohini Ramachandran. Localization of adaptive variants in human genomes using averaged
597 one-dependence estimation. *Nature communications*, 9(1):703, 2018.
- 598 The Anopheles gambiae 1000 Genomes Consortium. Ag1000G phase 1 AR3 data release. *Malar-*
599 *iaGEN*, 2015. URL <http://www.malariagen.net/data/ag1000g-phase1-AR3>.
- 600 Stefan Van Der Walt, S Chris Colbert, and Gael Varoquaux. The numpy array: a structure for
601 efficient numerical computation. *Computing in Science & Engineering*, 13(2):22, 2011.
- 602 Fernando A Villanea and Joshua G Schraiber. Multiple episodes of interbreeding between near-
603 derthal and modern humans. *Nature ecology & evolution*, 3(1):39, 2019.
- 604 Samuel K Wasser, Andrew M Shedlock, Kenine Comstock, Elaine A Ostrander, Benezeth Mutay-
605 oba, and Matthew Stephens. Assigning African elephant DNA to geographic region of origin:
606 applications to the ivory trade. *Proceedings of the National Academy of Sciences*, 101(41):
607 14847–14852, 2004.
- 608 S Wright. Isolation by distance. *Genetics*, 28(2):114–138, March 1943. URL [http://www.
609 genetics.org/cgi/reprint/28/2/114](http://www.genetics.org/cgi/reprint/28/2/114).
- 610 Sewall Wright. Isolation by distance under diverse systems of mating. *Genetics*, 31(1):336, 01
611 1946. URL <https://www.ncbi.nlm.nih.gov/pubmed/21009706>.
- 612 Wen-Yun Yang, John Novembre, Eleazar Eskin, and Eran Halperin. A model-based approach
613 for analysis of spatial structure in genetic data. *Nature Genetics*, 44:725 EP –, 05 2012. URL
614 <https://doi.org/10.1038/ng.2285>.
- 615 Tatiana Zerjal, Yali Xue, Giorgio Bertorelle, R Spencer Wells, Weidong Bao, Suling Zhu, Raheel
616 Qamar, Qasim Ayub, Aisha Mohyuddin, Songbin Fu, et al. The genetic legacy of the Mongols.
617 *The American Journal of Human Genetics*, 72(3):717–721, 2003.
- 618 Sha Joe Zhu, Jason A Hendry, Jacob Almagro-Garcia, Richard D Pearson, Roberto Amato, Alistair
619 Miles, Daniel J Weiss, Tim CD Lucas, Michele Nguyen, Peter W Gething, et al. The origins and
620 relatedness structure of mixed infections vary with local prevalence of *P. falciparum* malaria.
621 *eLife*, 8:e40845, 2019.
- 622 Veronika Zikmundová. *Spoken sibe: morphology of the inflected parts of speech*. Karolinum
623 Press, 2013. URL [https://www.google.com/books/edition/Spoken_Sibe_Morphology_of_
624 the_Inflected/PUs3BAAQBAJ?hl=en&gbpv=1](https://www.google.com/books/edition/Spoken_Sibe_Morphology_of_the_Inflected/PUs3BAAQBAJ?hl=en&gbpv=1).

625 Appendix 1: Theoretical limits on accuracy

626 Suppose that we know the spatial locations of some relatives of a given individual, and want to
627 predict the location of that focal individual. This is a best-case scenario for our actual problem, as
628 in fact we would have to infer the degrees of relatedness of the reference set to the focal individual,
629 but the calculations are useful in establishing a lower bound on the resolution of inference.

630 Suppose furthermore that the displacement in spatial position along each parent-child rela-
631 tionship has mean zero and variance σ , so that the net distance traveled along any path along k
632 links in the pedigree has mean zero and variance $k\sigma$. Given the location of n relatives of a focal
633 individual, a simple estimator of that individual's spatial location is simply the average of their
634 locations. How well does this do?

We can associate each link between parent p and child c in the pedigree with the displacement between them, $X_{pc} = -X_{cp}$; we have assumed that $\text{var}[X_{cp}] = \sigma^2$ for each. Suppose that the i^{th} relative can be reached by traversing relatives r_{i1}, \dots, r_{ik_i} , and so their location relative to the focal individual is $Y_i = X_{r_{i1}, r_{i2}} + \dots + X_{r_{i(k_i-1)}, r_{ik_i}}$. To compute the variance of our estimator, $\bar{Y} = \sum_{i=1}^n Y_i/n$, let n_{cp} be the number of i for which X_{cp} appears in the sum for Y_i , so that $\bar{Y} = \sum_{cp} n_{cp} X_{cp}/n$. Then, simply, $\text{var}[\bar{Y}] = \sum_{cp} (n_{cp}/n)^2 X_{cp}$. For instance, if those relatives are all 2^k ancestors k generations ago (i.e., the great $^{k-2}$ -grandparents) of the focal individual, then each of the 2^ℓ links between the ℓ^{th} and $(\ell - 1)^{\text{th}}$ generations are traversed by $2^{k-\ell}$ of the paths, and so

$$\text{var}[\bar{Y}] = \sum_{\ell=1}^k 2^\ell \left(\frac{2^{k-\ell}}{2^k} \right)^2 \sigma^2 = (1 - 2^{-k})\sigma^2.$$

635 Clearly, with less full pedigree coverage and more distant relatives, the error would become worse,
636 but it does not depend strongly on the degree of relatedness used: in general, using a few close or
637 many distant relatives should give an estimate of location within some moderate factor of σ .

638 **Supplementary Figures and Tables**

Neighborhood Size	Expected dispersal (map units/gen)	Error (map units) <i>mean (95% interval)</i>	Error (generations)
13	0.45	1.84 (0.259-3.633)	4.09
25	0.63	2.44 (0.388-6.033)	3.87
105	1.29	4.07 (0.685-8.874)	3.16
251	2.00	6.44 (0.639-14.526)	3.22
565	3.00	9.70 (0.871-21.146)	3.23

Table S1: Validation error in terms of map units and generations of dispersal for **Locator** runs in simulations with 450 training samples and 100,000 SNPs.

Species	kernel peak error (km) <i>median (90% interval)</i>	centroid error (km)
<i>Plasmodium falciparum</i>	16.92 (1.357 - 892.751)	218.98 (16.186 - 978.691)
<i>Anopheles gambiae/coluzzii</i>	5.69 (0.52 - 654.66)	36.03 (2.27 - 1579.79)
<i>Homo sapiens</i>	84.97 (4.42 - 2826.33)	452.62 (37.67 - 2178.94)

Table S2: Test error for windowed analyses of empirical datasets using the location with highest kernel density and the centroid of per-window predictions, as *median (90% interval)*.

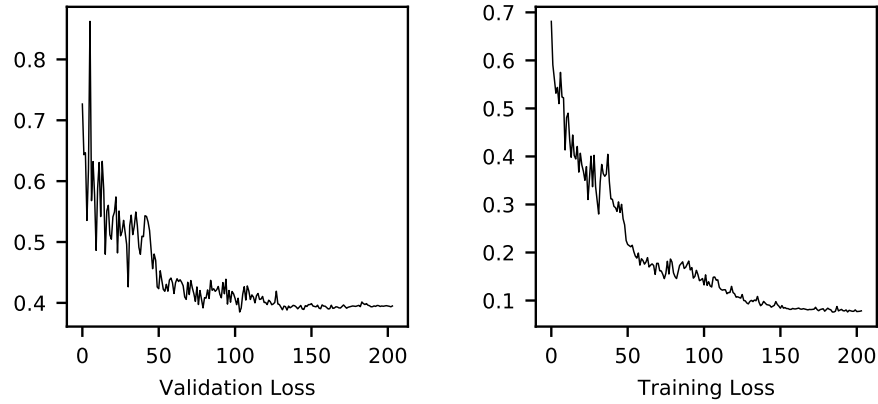


Figure S1: Example training and validation loss histories for a run on a single window of the dispersal=0.63 simulation. Epochs are shown on the horizontal axis and normalized loss on the vertical axis. The first three epochs (with very high loss) were excluded from the plot to improve axis scaling.

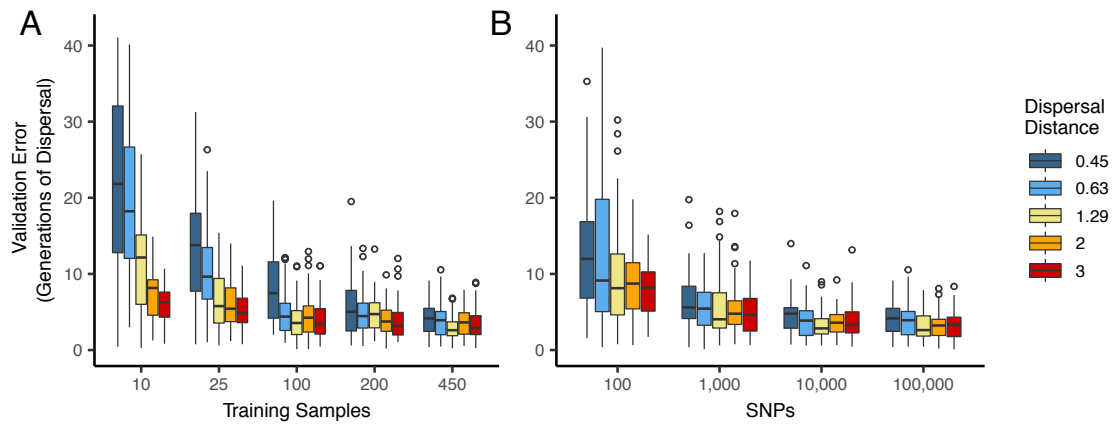


Figure S2: Validation error for *Locator* runs on simulations with varying dispersal distance, expressed in generations of mean dispersal (test error divided by mean dispersal distance per generation). A: Error for runs with 100,000 SNPs and varying numbers of training samples. B: Error for runs with 450 training samples and varying number of SNPs.

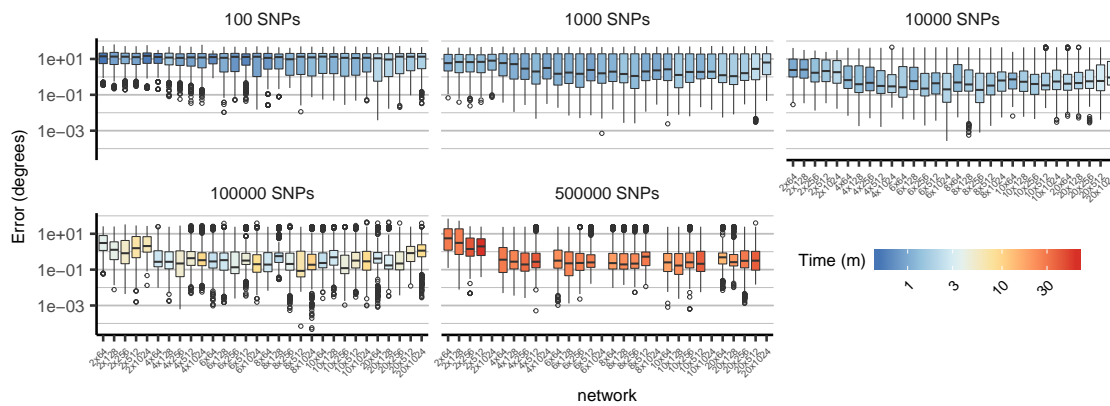


Figure S3: Comparison of cross-validation performance on the ag1000g dataset using SNPs from chromosome 3R, under varying network architectures and numbers of SNPs. Boxplots show the distribution of Euclidean distance between the true and predicted locations of validation samples across 10 replicate training runs. Network shapes are described on the horizontal axis as “layers × width”. Though 2-layer networks are typically the least accurate, no single architecture provides consistently better performance across datasets of different sizes. Missing networks required over 12GB GPU RAM.

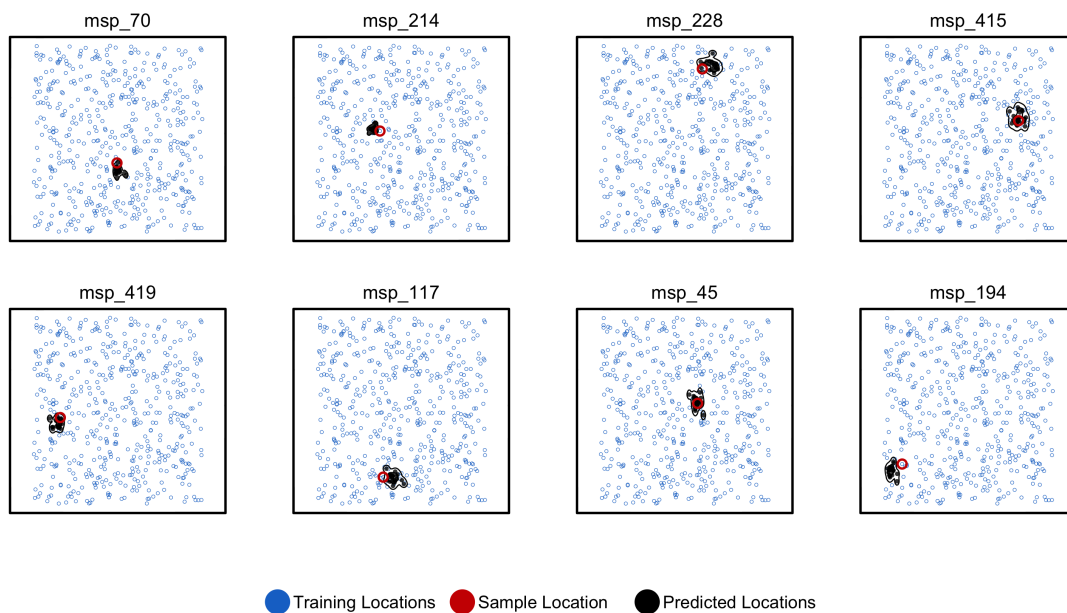


Figure S4: Predicted and true locations for 8 individuals simulated in a population with an expected dispersal rate of 0.63 map units / generation, using a set of 10,000 randomly sampled SNPs. Here we generate predictions (black points) from bootstrap samples of the complete genotype matrix (in contrast to using separate sets of SNPs extracted from windows as used for figures in the main text). This could be useful for low-density genotyping data from approaches like ddRADseq, or when users lack a reference genome for windowing. In this setting we see that the distribution of predictions is much smaller than fitting individual windows.

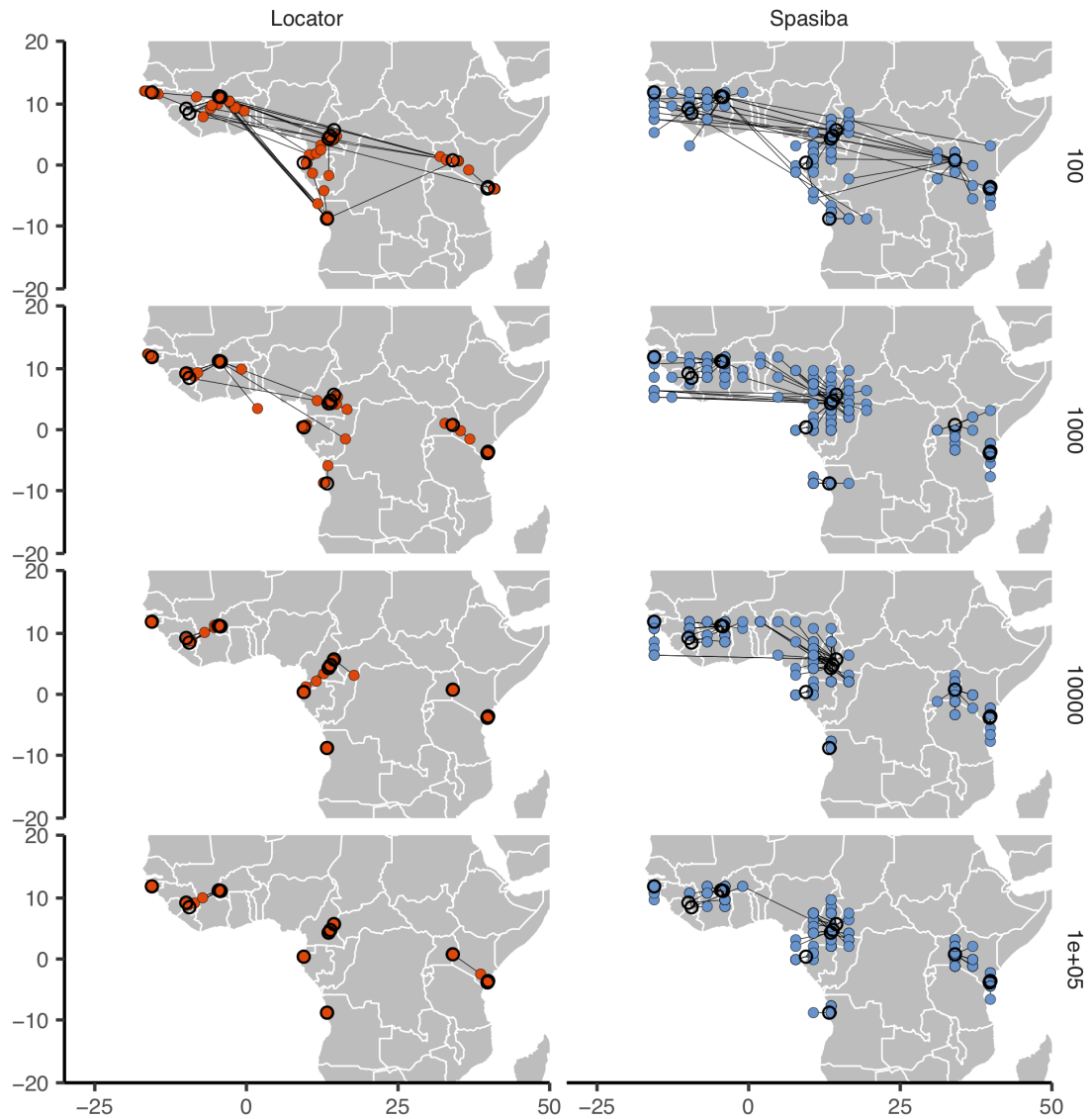


Figure S5: Predicted (colored points) and true (black circles) locations for Locator and SPASIBA on the ag1000g dataset. Number of SNPs per run is shown on the right. Both methods were run on randomly selected SNPs with minor allele count > 2 from the first five million base pairs of chromosome 2L.

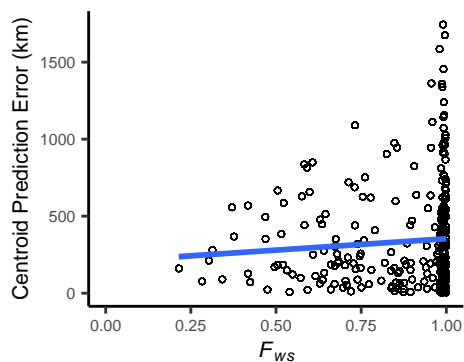


Figure S6: Centroid prediction error as a function of within-host diversity (F_{ws}) for the *Plasmodium falciparum* dataset. F_{ws} scales from 0 (maximum complexity) to 1 (minimum complexity). The blue line shows a linear regression ($p < 2.2e-16$, $R^2 = 0.006$, $slope = 148.1$). High within-host diversity does not appear to explain outliers in Locator's prediction error.

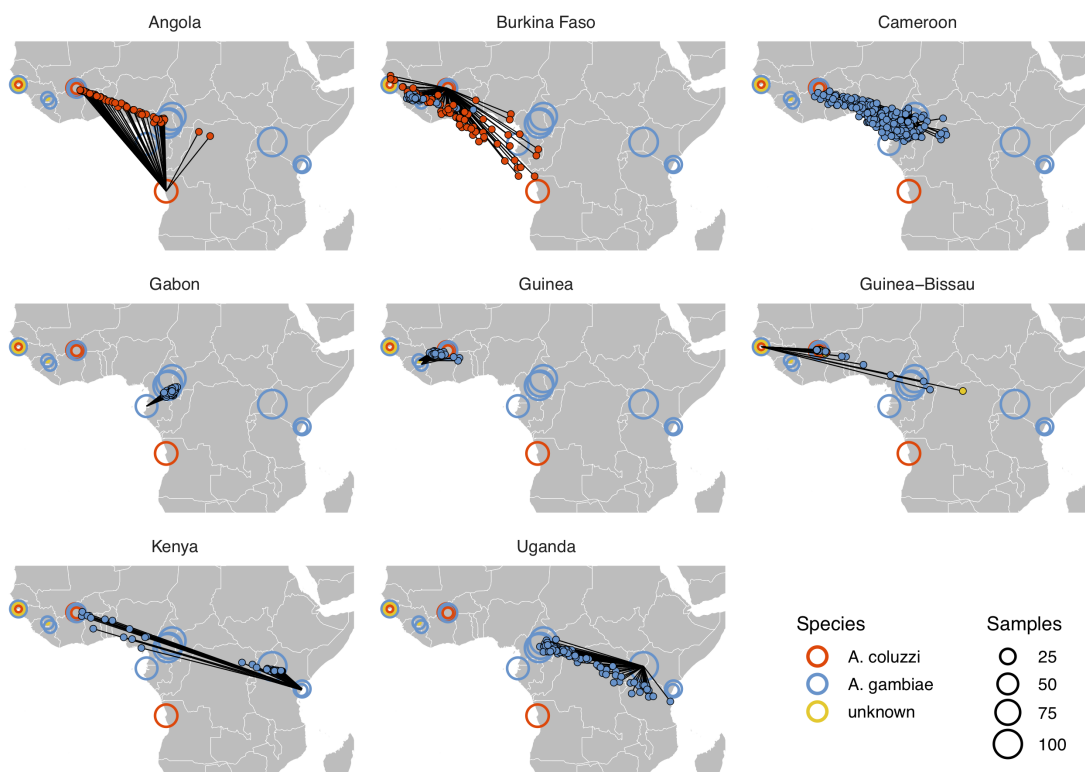


Figure S7: Performance on 10,000 SNPs from chromosome 2L in the ag1000g phase 1 dataset when all samples from localities in the true country are dropped from the training set.

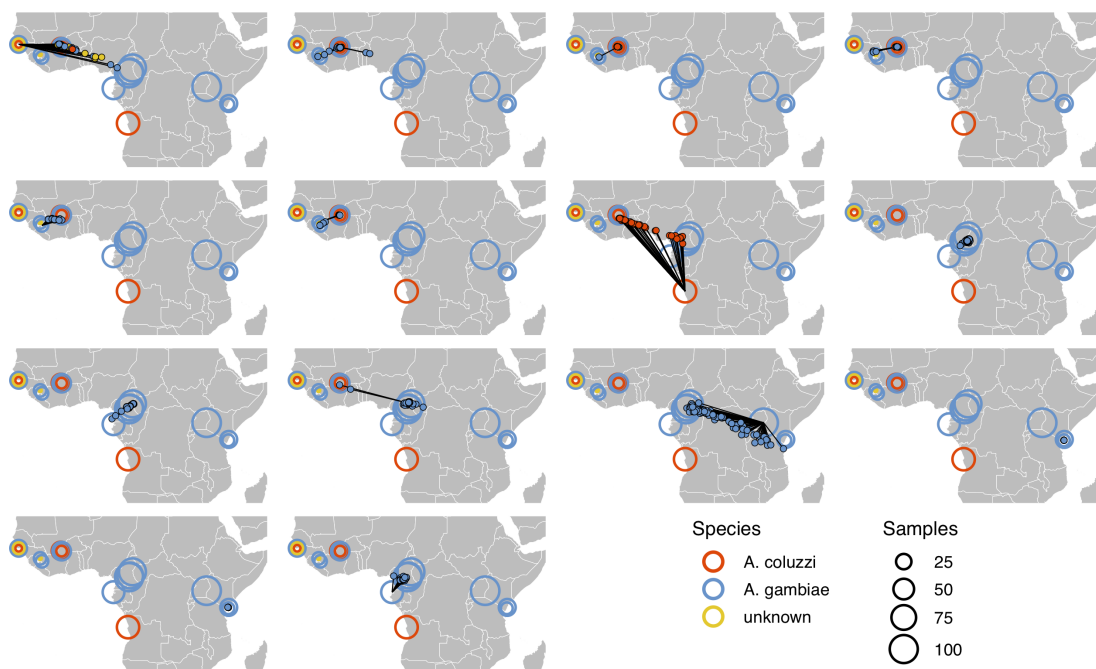


Figure S8: Performance on 10,000 SNPs from chromosome 2L in the ag1000g phase 1 dataset when all samples from the true locality are dropped from the training set.

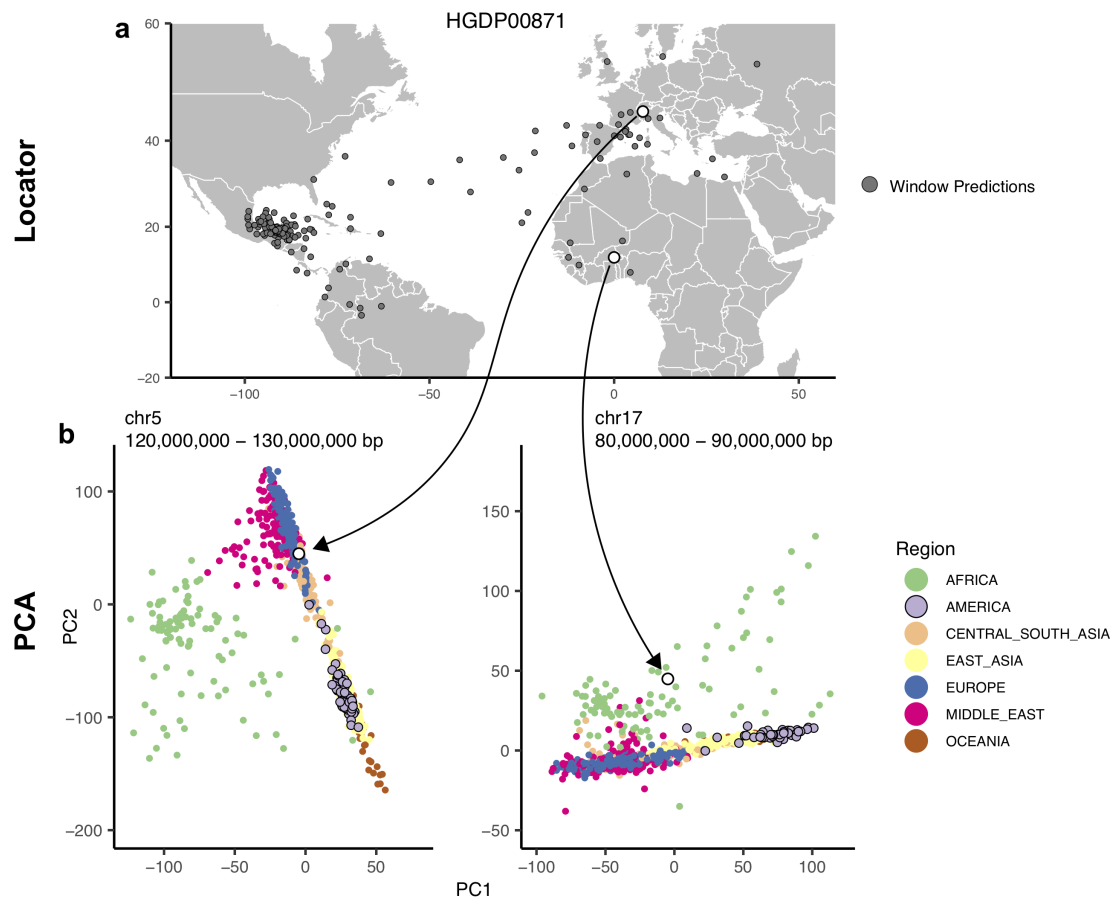


Figure S9: Outliers in windowed *Locator* analyses identify genomic regions enriched for admixed ancestry. A: Windowed *Locator* predictions for Maya sample HGDP00871. B: PCAs of all HGDP samples run on SNPs extracted from windows with predicted locations in western Europe (left) and west Africa (right). In these windows sample HGDP00871 (open points) clusters with individuals from region predicted by *Locator* in PC space, rather than with other genomes from the Americas.