

Visualizing Population Structure with Variational Autoencoders

C. J. Battey¹, Gabrielle C. Coffing¹, and Andrew D. Kern¹

¹University of Oregon Institute of Ecology and Evolution

October 20, 2020

Abstract

Dimensionality reduction is a common tool for visualization and inference of population structure from genotypes, but popular methods either return too many dimensions for easy plotting (PCA) or fail to preserve global geometry (t-SNE and UMAP). Here we explore the utility of variational autoencoders (VAEs) – generative machine learning models in which a pair of neural networks seek to first compress and then recreate the input data – for visualizing population genetic variation. VAEs incorporate non-linear relationships, allow users to define the dimensionality of the latent space, and in our tests preserve global geometry better than t-SNE and UMAP. Our implementation, which we call `popvae`, is available as a command-line python program at github.com/kr-colab/popvae. The approach yields latent embeddings that capture subtle aspects of population structure in humans and *Anopheles* mosquitoes, and can generate artificial genotypes characteristic of a given sample or population.

Introduction

As we trace the genealogy of a population forward in time, branching inherent in the genealogical process leads to hierarchical relationships among individuals that can be thought of as clades. Much of the genetic variation among individuals in a species thus reflects the history of isolation and migration of their ancestors. Describing this population structure is itself a major goal in biogeography, systematics, and human genetics; wherein one might attempt to infer the number of genotypic clusters supported by the data (Holsinger and Weir, 2009), estimate relative rates of migration (Petkova et al., 2016), or observe turnover in the ancestry of people living in a geographic region (Antonio et al., 2019).

Estimation of population structure is also critical for our ability to accurately link genetic variation to phenotypic variation, because population structure is a major confounding factor in genome-wide association studies (GWAS) (Lander and Schork, 1994; Pritchard and Donnelly, 2001; Marchini et al., 2004; Freedman et al., 2004). Downstream studies that use GWAS information can themselves be compromised by inadequate controls for structure, for instance in recent work trying to identify the effects of natural selection

34 on complex traits (Mathieson and McVean, 2012; Berg et al., 2019; Sohail et al., 2019).
35 Dimensionality reduction via principal components analysis (PCA) has been an important
36 tool for geneticists in this regard, and is now commonly used both to control for the effects
37 of population structure in GWAS (Price et al., 2006; Patterson et al., 2006) as well as for
38 visualization of genetic variation.

39 As a visualization tool however, PCA scatterplots can be difficult to interpret because
40 information about genetic variation is split across many axes, while efficient plotting is
41 restricted to two dimensions. Though techniques like plotting marginal distributions as
42 stacked density plots can aid interpretation, these require binning samples into "popula-
43 tions" prior to visualization, are rarely used in practice, and remain difficult to interpret
44 in complex cases. Recently two techniques from the machine learning community – t-SNE
45 (Maaten and Hinton, 2008) and UMAP (McInnes et al., 2018) – have shown promising per-
46 formance in producing two-dimensional visualizations of high-dimensional biological data.
47 In the case of UMAP, Diaz-Papkovich et al. (2019) recently showed that running the algo-
48 rithm on a large set of principal component axes allows visualization of subtle aspects of
49 population structure in three human genotyping datasets.

50 However, interpreting UMAP and t-SNE plots is also complicated by a lack of so-called
51 global structure. Though these methods perform well in clustering similar samples, distances
52 between groups are not always meaningful – two clusters separated by a large distance in a
53 t-SNE plot can be more similar to each other than either is to their immediate neighbors
54 (Becht et al., 2019). The degree to which initialization and hyperparameter tuning can
55 alleviate this issue remains an open question in the literature (Kobak and Linderman, 2019).

56 To create meaningful and interpretable visualizations of population genetic data we
57 would like a method that encodes as much information as possible into just two dimen-
58 sions while maintaining global structure. One way of achieving this is with a variational
59 autoencoder (VAE).

60 VAEs consist of a pair of deep neural networks in which the first network (the encoder)
61 encodes input data as a probability distribution in a latent space and the second (the de-
62 coder) seeks to recreate the input given a set of latent coordinates (Kingma and Welling,
63 2013). Thus a VAE has as its target the input data itself. The loss function for a VAE
64 is the sum of reconstruction error (how different the generated data is from the input)
65 and Kullback-Leibler (KL) divergence between a sample's distribution in latent space and
66 a reference distribution which acts as a prior on the latent space (here we use a standard
67 multivariate normal, but see (Davidson et al., 2018) for an alternative design with a hy-
68 perspherical latent space). The KL term of the loss function incentivizes the encoder to
69 generate latent distributions with meaningful distances among samples, while the recon-
70 struction error term helps to achieve good local clustering and data generation. VAE's have
71 been used extensively in image generation (Gulrajani et al., 2016; Larsen et al., 2015; Hou
72 et al., 2016) and several recent studies have applied them to dimensionality reduction and
73 classification of single-cell RNAseq data (Wang and Gu, 2018; Grønbech et al., 2018; Lafarge
74 et al., 2018; Hu and Greene, 2019). At deeper timescales than we test here, Derkarabetian
75 et al. (2019) recently explored the use of VAEs in species delimitation.

76 In population genetics two recent studies have studied the utility of generative deep
77 neural networks for creating simulated genotypes. Montserrat et al. (2019) use a class-
78 conditional VAE to generate artificial human genotypes, while Yelmen et al. (2019) use a
79 restricted Boltzman machine and provide an in-depth assessment of the population genetic
80 characteristics of their artificial genotypes. These studies found that such generative meth-
81 ods can produce short stretches of artificial genotypes that are difficult to distinguish from

82 real data, but performance was improved by using a generative adversarial network (GAN)
83 – either in combination with a VAE as in Montserrat et al. (2019) or as a standalone method
84 in Yelmen et al. (2019). In this study we focus not on generation of simulated genotypes,
85 but instead on the learned latent space representations of genotypes produced by a VAE,
86 and study when and how they can best be used for visualizing population structure.

87 We introduce a new method, `popvae` (for population VAE), a command-line python
88 program that takes as input a set of unphased genotypes and outputs sample coordinates
89 in a low-dimensional latent space. We test `popvae` with simulated data and demonstrate
90 its utility in empirical datasets of humans and *Anopheles* mosquitoes. In general `popvae` is
91 most useful for complex samples for which PCA projects important aspects of structure
92 across many axes. Relative to t-SNE and UMAP, the approach appears to better preserve
93 global geometry at the cost of less pronounced clustering of individual sample localities.
94 However, we show that hyperparameter tuning and stochasticity associated with train/test
95 splits and parameter initialization are ongoing challenges for a VAE-based method, and the
96 approach is much more computationally intensive than PCA.

97 Methods

98 Model

99 In this manuscript we describe the application of a Variational Auto-Encoder (VAE) to pop-
100 ulation genetic data for clustering and visualization Kingma and Welling (2013). Formally
101 let X be our dataset consisting of N observations (i.e. individual genotypes) such that
102 $X = \{x_1, x_2, \dots, x_N\}$, and let the probability of those data with some set of parameters θ be
103 $p_\theta(X)$. For VAEs we are interested in representing the data with a latent model, assigning
104 some latent process parameters z , such that we can write a generative latent process as
105 $p_\theta(x, z) = p_\theta(z)p_\theta(x|z)$, where $p_\theta(z)$ is the prior distribution on z . The last conditional
106 probability here $p_\theta(x|z)$ is often referred to as the decoder, as it maps from latent space to
107 data space.

108 For VAEs we also define a so-called encoder model $q_\phi(z|x)$, where ϕ represents the
109 parameters of the encoding (the mapping of x to the latent space z), and we seek to optimize
110 the encoder such that $q_\phi(z|x) \approx p_\theta(z|x)$. In practice the parameters ϕ represent the weights
111 and biases of the encoding neural network. We thus step from data space by using

$$(\boldsymbol{\mu}, \log(\boldsymbol{\sigma})) = \text{EncoderNeuralNetwork}(X) \quad (1)$$

$$q_\phi(z|x) = \mathcal{N}(z; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \quad (2)$$

112 The complete VAE information flow then has three steps: the encoder estimates sample
113 distributions in latent space as $q_\phi(z|x)$, we sample from the prior on the latent space using
114 $p_\theta(z)$, and finally decode back to data space using $p_\theta(x|z)$. Training is then performed by
115 optimizing the *evidence lower bound* or ELBO which has parameters of the encoder and
116 decoder within it such that

$$\mathcal{L}_{\theta, \phi}(X) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (3)$$

117 Optimization of the ELBO here leads to simultaneous fitting of the parameters of the en-
118 coder, ϕ , and the decoder, θ . In practice we use binary cross-entropy between true and
119 generated sequences for the first term, and Kullback-Leibler divergence of sample latent
120 distributions (relative to a standard normal $\mathcal{N}(0, 1)$) for the second term of equation 3. A
121 graphical depiction of this computational flow can be seen in Figure 1.

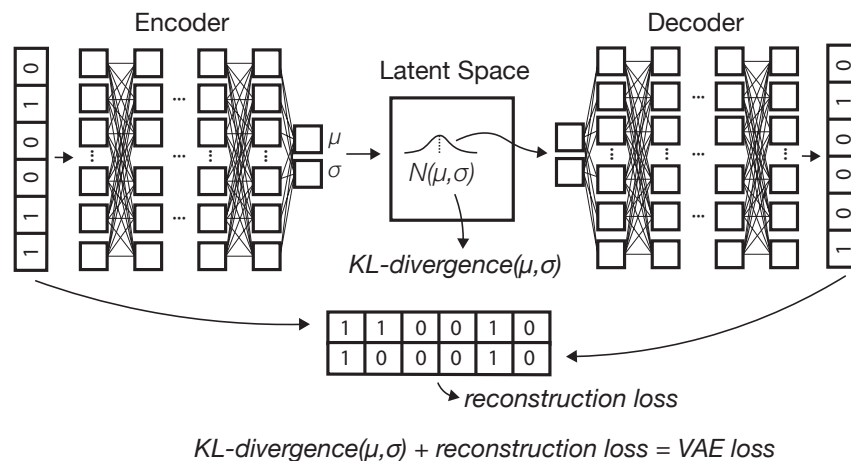


Figure 1: A schematic of the variational autoencoder (VAE) architecture. Input allele counts are passed to an encoder network which outputs parameters describing a sample’s location as a multivariate normal in latent space. Samples from this distribution are then passed to a decoder network which generates a new genotype vector. The loss function used to update weights and biases of both networks is the sum of reconstruction error (from comparing true and generated genotypes) and Kullback-Leibler divergence between sample latent distributions and $\mathcal{N}(0, 1)$.

122 Implementation

123 We implemented this model in python 3 using the tensorflow and keras libraries (Abadi et al.,
 124 2015; Chollet et al., 2015), with preprocessing relying on numpy, pandas, and scikit-allele
 125 (Miles and Harding, 2017; Oliphant, 2006–; McKinney, 2010). `popvae` reads in genotypes
 126 from VCFs, Zarr files <https://zarr.readthedocs.io/en/stable/>, or a bespoke hdf5 file
 127 format. Genotypes are first filtered to remove singletons and non-biallelic sites, and missing
 128 data is filled by taking two draws from a binomial distribution with probability equal to the
 129 allele frequency across all samples (a binned version of the common practice of filling missing
 130 genotypes with the mean allele frequency (Jombart, 2008; Dray and Josse, 2015)). Filtered
 131 genotypes are then encoded with 0/0.5/1 representing homozygous ancestral, heterozygous,
 132 and homozygous derived states, respectively.

133 Samples are split into training and validation sets before model training. We also ex-
 134 perimented with using all samples for training and a fixed number of epochs but found
 135 this generally led to poor performance (Appendix 1, Figure S1). Training samples are used
 136 to optimize weights and biases of the neural network, while validation samples are used
 137 to measure validation loss after each training epoch (a complete pass through the data),
 138 which in turn tunes hyperparameters of the optimizer. By default we use a random 90% of
 139 samples for training. However we found considerable variation in latent representations of
 140 some datasets when using different sets of training and validation samples (see e.g. Figure
 141 S2), so we encourage users to compare multiple training runs with different starting seeds

142 when interpreting plots.

143 `popvae`'s encoder and decoder networks are fully-connected feed-forward networks whose
144 size is controlled by two parameters – ‘width’, which sets the number of hidden units per
145 layer, and ‘depth’, which sets the number of hidden layers. We experimented with a range
146 of network sizes and set defaults to depth 6 and width 128, which performed well on the
147 empirical analyses described here (Table S1, Figure S3). However we also include a grid
148 search function by which `popvae` will conduct short training runs across a user-defined range
149 of network sizes and then fit a final model using the network size with minimum validation
150 loss.

151 We use a linear activation on the input layers to both networks and a sigmoid activation
152 on the output of the decoder (this produces numeric values bound by (0, 1)). We interpret the
153 sigmoid decoder outputs as the probability of observing a derived allele at a site, consistent
154 with our 0/0.5/1 encoding of the input genotypes. All other layers use “elu” activations
155 (Clevert et al., 2015), a modification of the more common “relu” activation which avoids
156 the “stuck neuron” problem by returning small but nonzero values with negative inputs.

157 We use the Adam optimizer (Kingma and Ba, 2014) and continue model training until
158 validation loss has not improved for p epochs, where p is a user-adjustable ‘patience’
159 parameter. We also set a learning rate scheduler to decrease the learning rate of the optimizer
160 by half when validation loss has not improved for $p/4$ epochs. This is intended to force the
161 optimizer to take small steps when close to the final solution, which increases training time
162 but in our experience leads to better fit models. Users can adjust many hyperparameters
163 from the command line, and modifying our network architectures is straightforward for those
164 familiar with the Keras library.

165 To evaluate model training `popvae` returns plots of training and validation loss by epoch
166 (e.g., Figure S4), and also outputs estimated latent coordinates for validation samples given
167 the encoder parameters at the end of each epoch. These can then be plotted to observe
168 how the model changes over the course of training, which can sometimes help to diagnose
169 overfitting. We also include an interactive plotting function which generates a scatter plot of
170 the latent space and allows users to mouse-over points to view metadata (Figure S5). This is
171 intended to allow users to quickly iterate through models while adjusting hyperparameters.
172 In Appendix 1 we discuss alternate approaches to network design and optimization tested
173 while developing `popvae`.

174 `popvae` is available at <https://github.com/kr-colab/popvae>, and scripts for re-
175 producing plots and analyses in this manuscript are available at [https://github.com/](https://github.com/cjbattey/popvae_analysis_scripts)
176 `cjbattey/popvae_analysis_scripts`. HGDP genotypes used in this paper are avail-
177 able at <ftp://ngs.sanger.ac.uk/production/hgdp>, AG1000G genotypes at [https://](https://www.malariagen.net/data/ag1000g-phase-2-ar1)
178 www.malariagen.net/data/ag1000g-phase-2-ar1, and 1000 genomes phase 3 data at
179 <https://www.internationalgenome.org/category/phase-3/>.

180 Results

181 Latent Spaces Reflect Human Migration History

182 We first applied `popvae` to 100,000 SNPs from chromosome 1 in the Human Genetic Diver-
183 sity Project (HGDP; Bergström et al. (2019)), a sample of global modern human diversity.
184 The resulting latent space reflects geography from the point of view of human demographic
185 history (Figure 2, Figure S6, Figure 4). Sub-Saharan African and South American popula-
186 tions are placed on opposite ends of one latent dimension, and north African (Mozabite) and

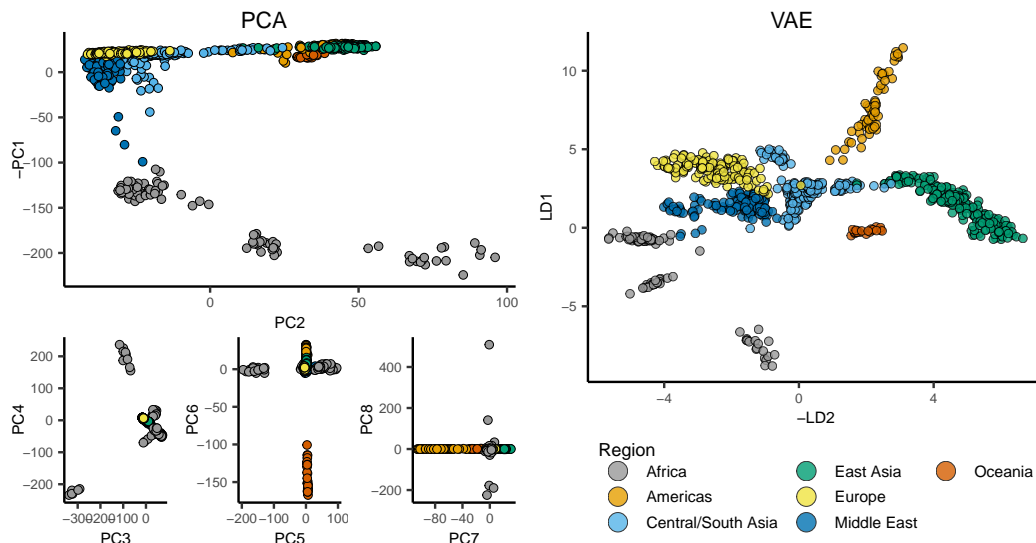


Figure 2: PCA axes 1-8 (left) and `popvae` run at default settings (right) for 100,000 random SNPs from chromosome 1 of the HGDP data. Axes are flipped to approximate geography.

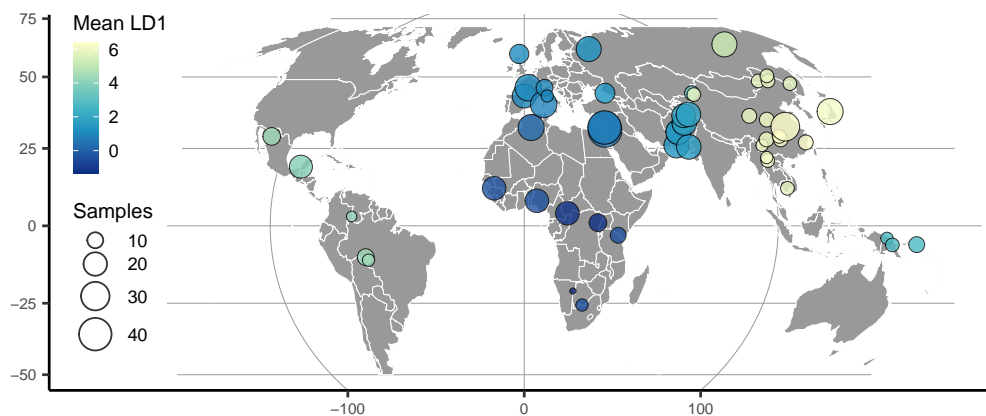


Figure 3: HGDP population locations with color scaled to the mean latent coordinate of a 1-dimensional `popvae` latent space.

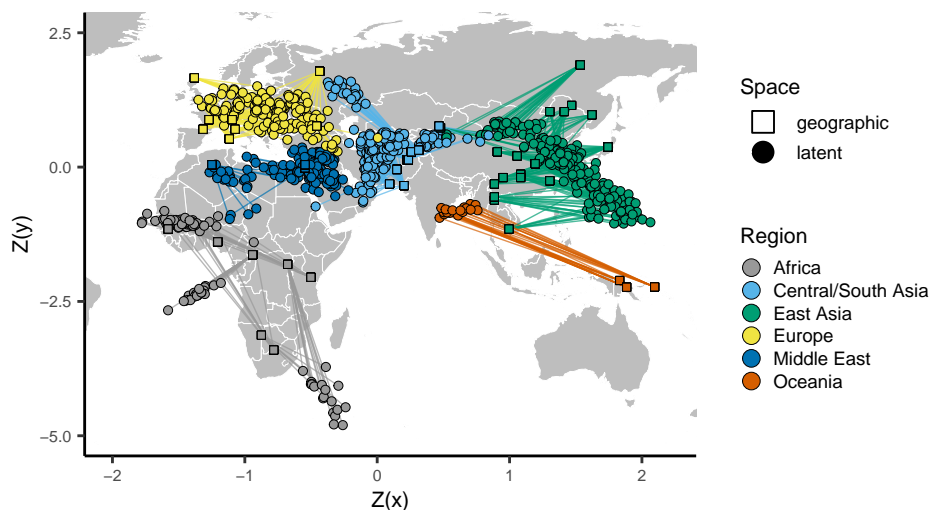


Figure 4: Comparing the VAE latent space with the geography of sampling localities in non-American HGDP samples (see Figure S8 for a plot including the Americas). Circles show z-normalized sample locations in latent space and squares show the corresponding location in geographic space.

187 east Asian samples are on opposite ends of the second; mirroring the geography of Africa
188 and Eurasia. Samples from the Americas are roughly centered among Eurasian samples on
189 latent dimension (LD) 2, consistent with recent demographic modeling studies suggesting
190 a mix of Eurasian ancestries in ancestral American populations (Flegontov et al., 2019;
191 Posth et al., 2018). Indeed the closest American samples to the European cluster are Maya
192 individuals who were found to have low levels of recent European admixture in previous
193 analyses (Bergström et al., 2019; Rosenberg et al., 2002) (Figure S6), suggesting *popvae*
194 is picking up on the signal of gene flow associated with European colonization of the Americas.

195 These patterns are similar to those seen in PCA, but many aspects of ancestry that are
196 difficult to see on the first two PC axes are conveniently summarized in *popvae*'s latent
197 space. For example, differentiation within the Americas and Oceania is not visible until
198 PC6 and PC7, respectively, but is clear in the 2D VAE latent space. This shows adjacent
199 clusters for the islands of Bougainville and Papua New Guinea, and a cline in Eurasian
200 ancestry from North through South America (Figure S6).

201 To highlight the flexibility of the VAE approach, we also trained a model with a 1-
202 dimensional latent space and used this to scale colors on a sampling map (Figure 3). This
203 results in a single latent dimension that approximates the diagonal of our 2D model, with
204 African and East Asian samples on either end of the spectrum. A comparison using PCA but
205 summarizing only the first principal component emphasizes diversity within Africa (Figure
206 S7) and provides little resolution for out-of-Africa groups.

207 Finally, to emphasize the correspondence of the VAE latent space with geography, we
208 can also directly compare geographic and latent spaces by rescaling both sets of coordinates
209 with a z-normalization and plotting them together on a map (Figure 4). As can be seen, the
210 visual correspondence between geographic and latent coordinates is striking in this case.

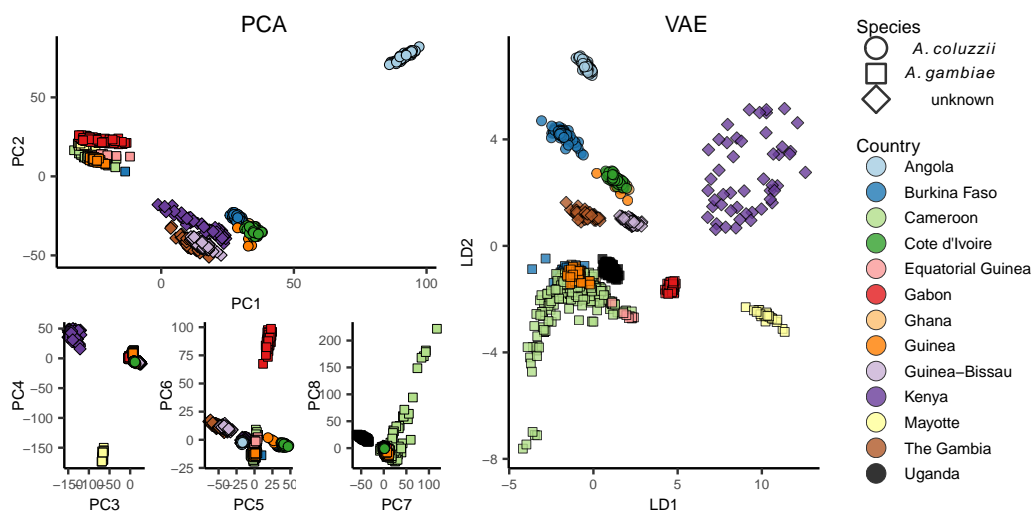


Figure 5: PCA (left) and VAE (right) run on 100,000 random SNPs from chromosome 3R of the AG1000G phase 2 data. Axes are flipped to approximate geography.

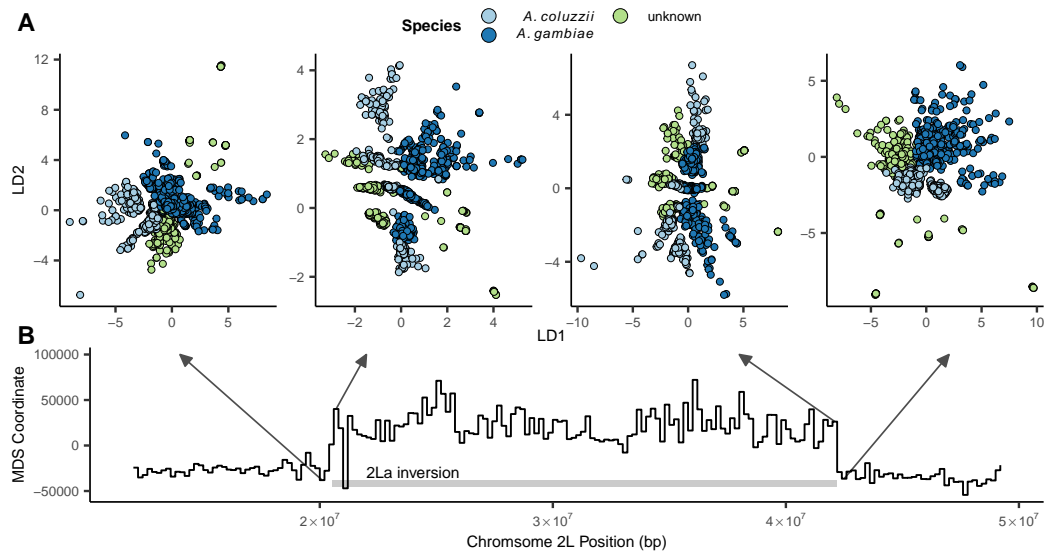


Figure 6: Latent spaces reflect inversion karyotypes at the 2La inversion in *A. gambiae* / *coluzzii*. A: VAE latent spaces for AG1000G phase 2 samples from windows near the 2La inversion breakpoints, colored by species. B: Multi-dimensional scaling values showing difference in the relative position of individuals in latent space across windows – high values reflect windows in which samples cluster by inversion karyotype, and low values by species.

211 Inversions and Population Structure in *Anopheles* Mosquitoes

212 We next applied `popvae` to DNA sequenced from the *Anopheles gambiae* / *coluzzii* complex
213 across sub-saharan African by the AG1000G project (AG1000G Consortium, 2020; Miles
214 et al., 2017) (Figure 5). Using 100,000 randomly-selected SNPs from chromosome 3R we
215 again find that the VAE captures elements of population structure that are not apparent
216 by visualizing two PC axes at a time. For example, samples from Kenya and the island of
217 Mayotte off East Africa are highly differentiated ($F_{st} > 0.18$ relative to all other groups),
218 but are placed between clusters of primarily west-African *coluzzii* and *gambiae* samples on a
219 plot of PC1/2. The VAE instead places these populations on the opposite end of one latent
220 dimension from all other groups and closest to Ugandan samples – similar to their relative
221 geographic position and positions on PC3/4. The VAE also captures the relatively high
222 differentiation of samples from Gabon and significant variation within Cameroon, which are
223 not visible until PC6 and PC8, respectively. Further details of population structure in this
224 species complex are discussed in AG1000G Consortium (2020).

225 *A. gambiae* / *coluzzii* genomes are characterized by a series of well-studied inversions on
226 chromosomes 2L and 2R (Coluzzi et al., 2002) which segregate within all populations and
227 are associated with both malaria susceptibility and ecological niche variation (Riehle et al.,
228 2017). The large 2La inversion contains at least one locus for insecticide resistance (*Rdl*),
229 and has experienced multiple hard sweeps and introgression events in its recent history
230 (Grau-Bové et al., 2020). Inversions have significant effects on local PCA (Li and Ralph,
231 2019) which often lead to samples clustering by inversion karyotype rather than geography
232 on the first two PC axes (Ma and Amos, 2012).

233 To test how our VAE responds to inversions we fit models to SNPs extracted from
234 200,000 bp non-overlapping windows across the 2LA inversion in the AG1000G phase 2 data
235 (Figure 6, Figure S11). We took an approach similar to Li and Ralph (2019) to summarize
236 differences in latent spaces across windows while accounting for axis rotation and scaling.
237 Latent dimensions were first scaled to 0 - 1 and the pairwise Euclidean distance matrix
238 among individuals was calculated for each window to generate rotation- and scale-invariant
239 representations of the latent space. We then calculated Euclidean distances among all pairs
240 of per-window distance matrices, giving us a matrix representing relative differences in latent
241 spaces across windows. Last, we used multi-dimensional scaling to compress this distance
242 matrix to a single dimension, and plotted this value against genomic position across the 2La
243 inversion region.

244 This analysis found two clear classes of latent spaces inside and outside the inversion
245 (Figure 6). Outside the inversion samples generally cluster by species and geography, while
246 inside the inversion samples form three clusters corresponding to the homozygous and het-
247 erozygous inversion karyotypes, similar to results found with PCA (Grau-Bové et al., 2020;
248 Riehle et al., 2017). Interestingly the VAE retains geographic and species clustering within
249 inversion classes, but loads these aspects of structure on a different latent dimension than
250 the karyotype clusters (e.g. LD1 reflects species clusters while LD2 reflects inversion karyo-
251 types in the windows shown in Figure 6). Unlike PCA, latent dimensions from a VAE
252 are not ranked by variance explained and nothing in the loss function incentivizes splitting
253 particular aspects of variation onto separate axes, so we found this pattern of partitioning
254 geographic and karyotypic signals somewhat surprising.

255 Simulations and Sensitivity Tests

256 In general a method's ability to detect population structure in a sample of genotypes scales
257 with the degree of differentiation and the size of the genotype matrix. Patterson et al. (2006)
258 found that there is a "phase change" phenomenon by which methods like PCA transition
259 from showing no evidence of structure to strong evidence of structure when $F_{st} \approx 1/\sqrt{nm}$,
260 where n is the number of genotyped SNPs and m is the number of sampled individuals.

261 To compare the performance of PCA and VAE around this threshold we ran a series of
262 two-population, isolation with migration model coalescent simulations in msprime (Kelleher
263 et al., 2016) while varying the symmetric migration rate to produce an expected equilibrium
264 F_{st} ranging from 0.0001 to 0.05. We sampled 50 diploid genomes from each population and
265 downsampled the resulting genotype matrix to 10,000 SNPs. Given this sample size we
266 expect the threshold for detecting structure to be approximately $F_{st} = 0.001$.

267 With tuned hyperparameters the VAE appeared slightly more sensitive to weak structure
268 than the first two axes of a PCA (Figure 7). Both `popvae` and PCA reflect some population
269 structure at $F_{st} \geq 0.005$ (though this is clearer in the VAE) but none at $F_{st} \leq 0.001$,
270 consistent with Patterson et al. (2006)'s "phase change" suggestion. However the VAE's
271 performance was highly sensitive to hyperparameter tuning on this dataset. At default
272 settings `popvae` latent spaces reflect no clear structure until $F_{st} = 0.05$ (Figure S12, Figure
273 S13). In particular we found that increasing the 'patience' parameter to 500 was necessary
274 for even marginal performance in this case, and running a grid search across network sizes
275 was needed to match PCA's sensitivity to weak structure.

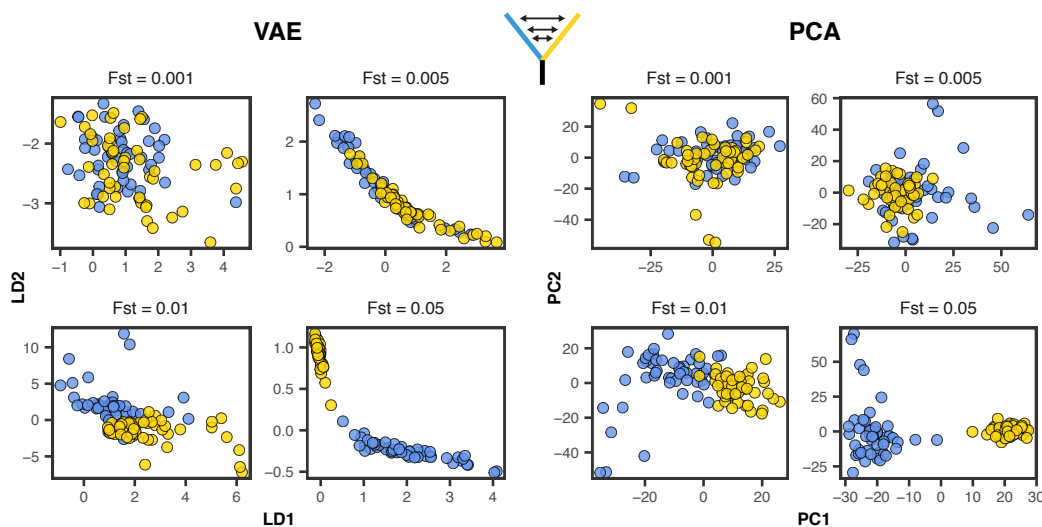


Figure 7: VAE latent spaces and PCA run on two-population coalescent simulations with F_{st} varying from 0.0001 – 0.05. Points are colored by population. `popvae` was run with tuned hyperparameters and patience set to 500. See Figure S12 for (much worse) performance with default settings.

276 Comparison with UMAP and t-SNE

277 In addition to PCA we also compared the VAE’s latent spaces to t-SNE (Maaten and Hinton,
278 2008) and UMAP (Diaz-Papkovich et al., 2019) (Figure S14, Figure S15), both of which
279 have been used recently for population genetic visualization. We first ran both methods on
280 the top 15 PC axes (following Diaz-Papkovich et al. (2019)) with default settings on the
281 human and *Anopheles* datasets and used the R packages ‘umap’ (Konopka, 2019) and ‘tsne’
282 (Donaldson, 2016) as our reference implementations.

283 For HGDP data both UMAP and t-SNE produce latent spaces that roughly correspond to
284 continental regions (Figure S14). Running both methods at default settings, UMAP’s latent
285 space was much more tightly clustered – for example grouping all samples from Africa into
286 a single small region. Similar patterns were seen in the AG1000G data (Figure S15) – both
287 t-SNE and UMAP produce latent spaces that strongly cluster sample localities and species.
288 However, global geometry appeared to be poorly preserved in t-SNE and UMAP latent
289 spaces. That is, though clusters in latent space correspond to sampling localities, distances
290 among clusters do not appear to meaningfully reflect geography or genetic differentiation.

291 To compare how well different methods reflect geography we compared pairwise distances
292 among individuals in latent and geographic space for Eurasian human samples (HGDP
293 regions Europe, Central/South Asia, the Middle East, and East Asia). Geographic distances
294 were great-circle distance calculated on a WGS84 ellipse with the R package ‘sp’ (Pebesma
295 et al., 2012). Distances were scaled to 0-1 for this analysis, and we calculated the coefficient
296 of determination (R^2) across geographic and latent-space distance for each method as a
297 metric. VAE latent space distances have the strongest correlation with geographic distance
298 (Figure 8; $R^2 = 0.659$), followed by PCA ($R^2 = 0.561$), UMAP ($R^2 = 0.529$), and t-SNE
299 ($R^2 = 0.342$).

300 Finally to test how parameter tuning of tSNE and UMAP impacts our results, we repro-
301 duced our analysis of HGDP data using double and triple the default values for `n_neighbors`
302 (UMAP) and `perplexity` (tSNE). Though scatter plots are visually similar at these settings
303 (Figure S16) the correlation between latent-space and geographic distances of Eurasian sam-
304 ples is improved in both methods at double default settings (t-SNE: $R^2 = 0.631$, UMAP:
305 $R^2 = 0.611$; Figure S17). At triple default settings we observed slightly better performance
306 for tSNE and slightly worse for UMAP (Figure S18, Figure S19).

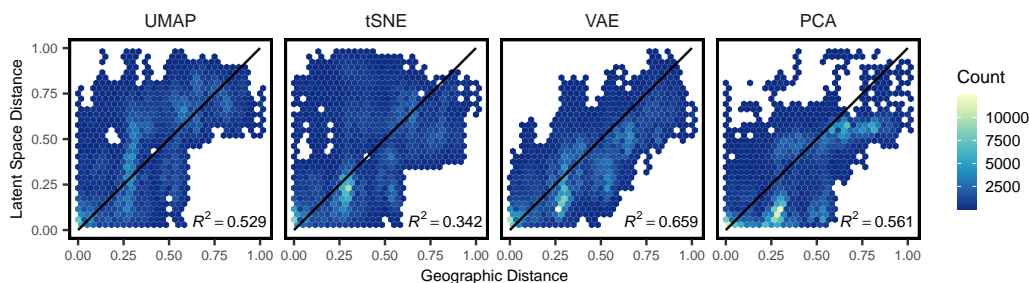


Figure 8: Comparing pairwise distances in geographic and latent space for Eurasian human genotypes across four dimensionality reduction methods run at default settings. All distances are scaled to 0-1. Black lines show a 1:1 relationship.

307 Run Times and Computational Resources

308 We compared `popvae`'s run times to PCA, UMAP, and t-SNE using sets of 100,000 and
309 10,000 SNPs from the HGDP as described above. `popvae` was run using default settings
310 (i.e. fitting a single network rather than running a grid search over network sizes) using
311 a consumer GPU (Nvidia GeForce RTX 2070). PCAs were run in the python package
312 `scikit-allel` (Miles and Harding, 2017), which in turn relies on singular-value decomposition
313 functions from the `numpy` library (Oliphant, 2006–).

314 `popvae` was much slower than PCA or UMAP, but comparable to running t-SNE on
315 PC coordinates. However for datasets of the size we tested here none of these run times
316 present significant challenges – all methods return sample latent coordinates in less than
317 five minutes. We have not conducted exhaustive tests on CPU training times for `popvae`,
318 but in general find these to require at least twice as much time as GPU runs.

319 However for larger datasets we expect `popvae`'s run time performance would suffer fur-
320 ther in comparison to PCA and UMAP. The major computational bottleneck is loading
321 tensors holding weights for the input and output layers of the encoder and decoder net-
322 works into GPU memory. These tensors have dimensions `n_snps x network_width` so they
323 become extremely large when running on large genotype matrices. Our development ma-
324 chine has 8GB GPU RAM and can process up to roughly 700,000 SNPs in a single analysis
325 using a 128-unit-wide network. Throughout this study we have limited our analysis to rela-
326 tively small subsets of genome-wide SNPs to allow us to explore a range of network sizes in
327 reasonable time. Scaling up to a single model fit to all genome-wide SNPs – on the order
328 of 10^7 for datasets like the HGDP – would require access to specialized hardware with very
329 large GPU memory pools.

run time (s)	SNPs	method
204.4	100,000	VAE
3.6		PCA
6		UMAP
124.8		t-SNE
78.8	10,000	VAE
0.5		PCA
2.7		UMAP
119.5		t-SNE

Table 1: Run times for VAE, PCA, UMAP, and t-SNE HGDP data. UMAP and t-SNE were run on the top 20 PC axes (run times thus include running the PCA).

330 Generating Genotypes

331 The VAE framework also allows us to generate genotypes characteristic of a given population
332 by sampling from the latent space of a trained model. Simulated genotypes generated by
333 process-based models like the coalescent are a key tool in population genetics, because they
334 allow us to explore the impact of various generative processes – demography, selection, etc –
335 on observed genetic variation (Adrión et al., 2020a). In contrast `popvae`'s generative model
336 provides essentially no mechanistic insight beyond the strong observed correlation of latent
337 and geographic spaces. However, if the VAE accurately reproduces characteristics of real

338 genotypes it could be a fast alternative to simulation that does not require parameterizing
339 a custom demographic model.

340 We compared these approaches by analyzing empirical data from European (CEU), Han
341 (CHB), and Yoruban (YRI) human genotypes in the 1000 Genomes Project data (Con-
342 sortium et al., 2015). We first subset 50 samples from each population and then fit a
343 2-dimensional `popvae` model to all SNPs from chromosome 22. To generate genotypes we
344 drew a sample from the latent distribution of each individual and passed these coordinates
345 to the trained decoder network. We interpret the sigmoid activation output of our decoder
346 as the probability of observing a derived allele at each site, and generate derived allele counts
347 by taking two draws from a binomial distribution with $p = g_{i,j}$ where $g_{i,j}$ is the decoder
348 output for individual i at site j .

349 As a baseline comparison we used coalescent simulations from the `standardpopsim` li-
350 brary (Adrión et al., 2020a) of the 3-population out-of-Africa model (`OutOfAfrica_3G09`) –
351 a rigorously tested implementation of the demographic model fit to the joint site frequency
352 spectrum in Gutenkunst et al. (2009) using the `msprime` coalescent simulator (Kelleher
353 et al., 2016). For this comparison we changed `standardpopsim`'s default human mutation
354 rate of 1.29×10^{-8} to 2.35×10^{-8} to match the rate used in Gutenkunst et al. (2009), used
355 the `HapMapII_GRCh37` recombination map for chromosome 22, and sampled 100 haploid
356 chromosomes from each population.

357 Last, we examined three facets of population genetic variation on real, VAE-generated,
358 and simulated genotype matrices: the site frequency spectrum, the decay of linkage disequi-
359 librium with distance along the chromosome, and embeddings from a PCA. These analyses
360 were conducted in `scikit-allel` (Miles and Harding, 2017) after masking genotypes to retain
361 only sites with the most stringent site accessibility filter ("P") in the 1000 genome project's
362 phase 3 site accessibility masks. LD statistics were calculated only for YRI samples using
363 SNPs between positions 2.5×10^7 and 2.6×10^7 in the hg18 reference genome and summa-
364 rized by calculating the mean LD for all pairs of alleles in 25 distance bins (similar results
365 in three different genomic windows are shown in figure S20). Results are plotted in figure 9.

366 In general we found all methods produce similar results in a plot of the first two PC axes,
367 suggesting they capture broad patterns of allele frequency variation created by population
368 structure. The site frequency spectrum is also very similar for the VAE and real data, while
369 the simulated genotypes suffer from a scaling issue. This could reflect differences in the
370 input data – Gutenkunst et al. (2009) fit models to an SFS calculated from a set of sanger-
371 sequenced loci in 1000 genomes samples, rather than the short-read resequenced SNPs from
372 the 1000 Genomes project we use – or an inaccuracy in one of the constants used to convert
373 scaled demographic model parameters to real values (accessible genome size, generation
374 time, or mutation rate). LD decay shows the largest difference among methods. Simulation
375 and real data both reflect higher LD among nearby SNPs which decays with distance, while
376 the VAE genotypes produced no correlation between distance along a chromosome and
377 pairwise LD.

378 These differences reflect the strengths and weaknesses of each method. The VAE decoder
379 doesn't require a pre-defined demographic model and by design exactly fits the matrix size
380 of input empirical data, so it should not suffer from the scaling issues that frequently impact
381 population genetic models. But the lack of mechanistic biological knowledge in its design
382 means it misses obvious and important features of real sequence data like the decline of
383 LD with distance. In this case the lack of LD decay in VAE decoder sequences means this
384 implementation should not be used for testing properties of analyses like GWAS, in which
385 LD among a subset of sequenced loci and an unknown number of truly causal loci is a crucial

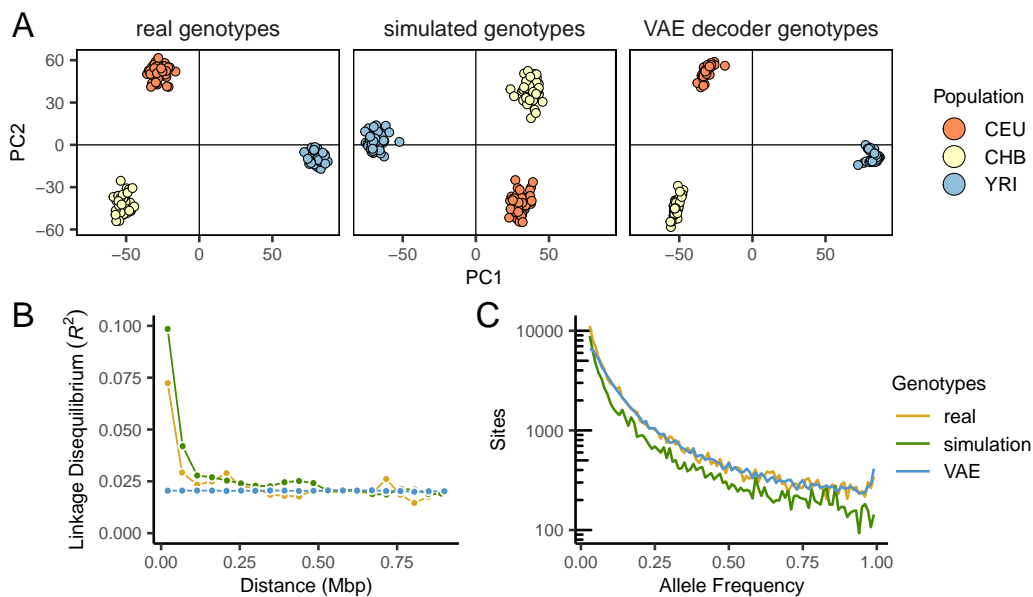


Figure 9: Comparing real, VAE-generated, and simulated genotype matrices for three populations from the 1000 genomes project. The VAE decoder and coalescent simulation produce similar results in genotype PCA (A), but the VAE fails to reproduce the decay of LD with distance along the chromosome seen in real data (B). The site frequency spectrum is very similar for real and VAE-generated genotypes, but suffers from scaling issues in the coalescent simulation (C).

386 parameter. Though other network designs (e.g. a convolutional neural network Fligel et al.
387 (2019) or a recurrent neural network Adrion et al. (2020b)) could potentially address the
388 specific shortcoming of LD decay, the general problem of a non-mechanistic generator failing
389 to mimic features of the data produced by well-understood processes seems intrinsic to the
390 machine learning approach.

391 Discussion

392 Dimensionality reduction of genotypic variation is a key analytic tool in modern genomics
393 and their visualizations are often the central figure of a genetic study. For example, Antonio
394 et al. (2019) studied a 10,000-year transect of genotypes from Rome and extensively used
395 PCA to visualize changes in ancestry in the city over time. In cases like this producing
396 informative plots of population structure is a requisite step for the analysis and can shape
397 the way data is interpreted both by authors and readers.

398 In this study we demonstrate how variational autoencoders can be used for visualization
399 and low dimensional summaries of genotype data. Variational autoencoders have at least
400 two attractive properties for genetic data: they allow users to define the output dimension-
401 ality, and they preserve global geometry (i.e., relative positions in latent space) better than
402 competing methods. As we have shown in humans and mosquitoes, this allows users to gen-
403 erate visualizations that summarize relationships among samples without either comparing
404 across several panels (as with PCA) or attempting to ignore possibly spurious patterns of
405 global structure (as with t-SNE and UMAP).

406 An additional attractive property of VAEs is that they are generative models. That is
407 to say that VAEs allow us to create genotypes that capture aspects of population genetic
408 variation characteristic of the training set. This is done by taking samples from the estimated
409 latent space and passing forward into data space. Though in theory this could be used as an
410 alternative to simulation, our implementation fails to replicate at least one important aspect
411 of real genomes – the decay of linkage disequilibrium with distance along a chromosome – and
412 thus offers limited utility for tasks such as boosting GWAS sample sizes or as a substitute
413 for simulation. We point researchers interested in generating genotypes via deep learning
414 approaches to recent work by Yelmen et al. (2019) and Montserrat et al. (2019), which
415 describe similar, deep learning based methods more tightly aimed at generating realistic
416 genotypes.

417 There are also several significant limitations of our method as a visualization tool. One
418 issue is that we lack a principled understanding of how the VAE output maps to parameters
419 of idealized population models like the coalescent (Kingman, 1982). This is in contrast to
420 PCA, which was first applied to genetic data with little theoretical background (Menozzi
421 et al., 1978) but is now fairly well characterized in reference to population genetic models
422 (McVean, 2009; Novembre and Stephens, 2008).

423 Hyperparameter tuning is another challenge. As we showed, `popvae` has many hyperpa-
424 rameters that significantly affect the output latent space and no principled way to set them
425 *a priori*. Though we include a grid-search function for network sizes, this is slow and is
426 still dependent on other hyperparameters – like the patience used for early stopping, or the
427 learning rate of the optimizer – which we have set to defaults that may not be optimal for
428 all datasets. This is not a unique issue to VAEs; opaque hyperparameters of methods like
429 t-SNE and UMAP can significantly affect embeddings (Kobak and Linderman, 2019), and
430 preprocessing choices such as how to scale allele counts prior to PCA dramatically vary the

431 appearance of final plots (Patterson et al., 2006). However it does require extra work on
432 the part of users interested in exploring the full parameter space.

433 A parallel issue is stochasticity in the output. Stochasticity is introduced by the random
434 test/train split, parameter initialization states, and even the execution order of operations
435 run in parallel on GPU during model training. Though all but the last of these can be fixed
436 by setting a random seed, which itself could be (unfortunately) seen as a hyperparameter,
437 there is no obvious way to compare models fit to different validation sets in a world of limited
438 training examples. This introduces noise which could potentially allow users to cherry-pick
439 a preferred latent space.

440 For example, one run of our best-performing network architecture on the HGDP data
441 produced a latent space with in which samples Papua New Guinea and Bougainville are
442 separated by roughly the same distance as samples from north Africa and East Asia. In
443 contrast all other fits of the same network architecture cluster these samples (Figure S2,
444 see the top middle panel). We chose a latent space for the main text that lacked this
445 feature because it occurred in only one training run, but acknowledge this procedure is
446 sub-optimal. Developing a method to summarize across multiple latent spaces, perhaps via
447 ensemble learning approaches, would be useful for postprocessing VAE output when latent
448 spaces vary.

449 The last major shortcoming is computational effort. `popvae` is much slower and more
450 computationally intensive than PCA, and requires specialized and expensive GPU or TPU
451 hardware to run on large sets of SNPs. Future developments in both hardware and software
452 will likely alleviate this issue somewhat, but at present it may make the method difficult
453 to apply to the increasingly common whole genome resequencing data now being generated
454 for many species.

455 One important question we did not explore in this study is whether VAE latent space co-
456 ordinates offer any improvement over PCA when used as covariates to correct for population
457 structure in GWAS (Price et al., 2006). UMAP and t-SNE are generally thought to be inap-
458 propriate for this use because of their failure to preserve global geometry (Diaz-Papkovich
459 et al., 2019), but because the VAE appears to strongly reflect geography in humans it may
460 be useful for this task. Testing this aspect of the VAE could be done in simulation but would
461 benefit from empirical investigations in large human datasets – a task which is beyond the
462 scope of the present study, but perhaps fruitful for further investigation.

463 Here we have shown that our implementation of a VAE, `popvae`, can produce informative
464 visualizations of population genetic variation and offers some benefits relative to competing
465 methods. However our approach is just one implementation of a huge class of potential
466 models falling under the VAE umbrella. Altering the prior on the latent space (Davidson
467 et al., 2018), the weighting of the loss function (Higgins et al., 2017), or the type of neural
468 network used in either the encoder or decoder all offer avenues for further research and
469 potential improvement (see also Appendix 1, where we briefly describe alternate approaches
470 we experimented with). Entirely different methods of visualizing population structure which
471 focus on genetic variants rather than individuals, like that proposed in Biddanda et al.
472 (2020), also offer a complementary perspective on the nature of genetic differentiation.
473 As population genetic data becomes increasingly common across evolutionary biology we
474 anticipate visualization techniques will receive increased attention from researchers in many
475 areas, and believe VAEs offer a promising avenue for research.

⁴⁷⁶ Acknowledgements

⁴⁷⁷ We thank Peter Ralph for the suggestion to use binomial sampling to bin the decoder
⁴⁷⁸ output, and other members of the Kern-Ralph co-lab for comments on the software and
⁴⁷⁹ manuscript. Thanks to the reviewers and editors for your comments. Audry Gill developed
⁴⁸⁰ a successful pilot of this project early after its inception. CJB was supported by NIH awards
⁴⁸¹ R01GM117241 and F32GM136123. ADK was supported under NIH awards R01GM117241
⁴⁸² and R01HG010774.

REFERENCES

REFERENCES

483 **References**

- 484 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
485 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian
486 Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz,
487 Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry
488 Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya
489 Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda
490 Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and
491 Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems,
492 2015. URL <http://tensorflow.org/>. Software available from tensorflow.org.
- 493 Jeffrey R Adrion, Christopher B Cole, Noah Dukler, Jared G Galloway, Ariella L Gladstein,
494 Graham Gower, Christopher C Kyriazis, Aaron P Ragsdale, Georgia Tsambos, Franz
495 Baumdicker, et al. A community-maintained standard library of population genetic mod-
496 els. *BioRxiv*, pages 2019–12, 2020a.
- 497 Jeffrey R Adrion, Jared G Galloway, and Andrew D Kern. Predicting the landscape of
498 recombination using deep learning. *Molecular biology and evolution*, 37(6):1790–1808,
499 2020b.
- 500 AG1000G Consortium. Genome variation and population structure among 1142 mosquitoes
501 of the african malaria vector species *Anopheles gambiae* and *Anopheles coluzzii*. *Genome*
502 *Research*, 2020. doi: 10.1101/gr.262790.120. URL [http://genome.cshlp.org/content/
503 early/2020/09/25/gr.262790.120.abstract](http://genome.cshlp.org/content/early/2020/09/25/gr.262790.120.abstract).
- 504 Margaret L Antonio, Ziyue Gao, Hannah M Moots, Michaela Lucci, Francesca Can-
505 dilio, Susanna Sawyer, Victoria Oberreiter, Diego Calderon, Katharina Devitofranceschi,
506 Rachael C Aikens, et al. Ancient rome: A genetic crossroads of europe and the mediter-
507 ranean. *Science*, 366(6466):708–714, 2019.
- 508 Etienne Becht, Leland McInnes, John Healy, Charles-Antoine Dutertre, Immanuel WH
509 Kwok, Lai Guan Ng, Florent Gehrmann, and Evan W Newell. Dimensionality reduction for
510 visualizing single-cell data using umap. *Nature biotechnology*, 37(1):38, 2019.
- 511 Jeremy J Berg, Arbel Harpak, Nasa Sinnott-Armstrong, Anja Moltke Joergensen,
512 Hakhamanesh Mostafavi, Yair Field, Evan August Boyle, Xinjun Zhang, Fernando
513 Racimo, Jonathan K Pritchard, et al. Reduced signal for polygenic adaptation of height
514 in uk biobank. *ELife*, 8:e39725, 2019.
- 515 Anders Bergström, Shane A. McCarthy, Ruoyun Hui, Mohamed A. Almarri, Qasim Ayub,
516 Petr Danecek, Yuan Chen, Sabine Felkel, Pille Hallast, Jack Kamm, Hélène Blanché, Jean-
517 François Deleuze, Howard Cann, Swapnil Mallick, David Reich, Manjinder S. Sandhu,
518 Pontus Skoglund, Aylwyn Scally, Yali Xue, Richard Durbin, and Chris Tyler-Smith. In-
519 sights into human genetic variation and population history from 929 diverse genomes.
520 *bioRxiv*, 2019. doi: 10.1101/674986. URL [https://www.biorxiv.org/content/early/
521 2019/06/27/674986](https://www.biorxiv.org/content/early/2019/06/27/674986).
- 522 Arjun Biddanda, Daniel P Rice, and John Novembre. Geographic patterns of human allele
523 frequency variation: a variant-centric perspective. *BioRxiv*, 2020.
- 524 François Chollet et al. Keras. <https://github.com/fchollet/keras>, 2015.

REFERENCES

REFERENCES

- 525 Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep
526 network learning by exponential linear units (ELUs). *arXiv preprint arXiv:1511.07289*,
527 2015.
- 528 Mario Coluzzi, Adriana Sabatini, Alessandra della Torre, Maria Angela Di Deco, and Vin-
529 cenzo Petrarca. A polytene chromosome analysis of the anopheles gambiae species com-
530 plex. *Science*, 298(5597):1415–1418, 2002.
- 531 1000 Genomes Project Consortium et al. A global reference for human genetic variation.
532 *Nature*, 526(7571):68–74, 2015.
- 533 Tim R Davidson, Luca Falorsi, Nicola De Cao, Thomas Kipf, and Jakub M Tomczak.
534 Hyperspherical variational auto-encoders. *arXiv preprint arXiv:1804.00891*, 2018.
- 535 Shahan Derkarabetian, Stephanie Castillo, Peter K Koo, Sergey Ovchinnikov, and Mar-
536 shal Hedin. A demonstration of unsupervised machine learning in species delimitation.
537 *Molecular phylogenetics and evolution*, 139:106562, 2019.
- 538 Alex Diaz-Papkovich, Luke Anderson-Trocmé, and Simon Gravel. Umap reveals cryptic pop-
539 ulation structure and phenotype heterogeneity in large genomic cohorts. *PLoS genetics*,
540 15(11), 2019.
- 541 Justin Donaldson. *tsne: T-Distributed Stochastic Neighbor Embedding for R (t-SNE)*, 2016.
542 URL <https://CRAN.R-project.org/package=tsne>. R package version 0.1-3.
- 543 Stéphane Dray and Julie Josse. Principal component analysis with missing values: a com-
544 parative survey of methods. *Plant Ecology*, 216(5):657–667, 2015.
- 545 Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The unreasonable effectiveness of
546 convolutional neural networks in population genetic inference. *Molecular biology and*
547 *evolution*, 36(2):220–238, 2019.
- 548 Pavel Flegontov, N Ezgi Altınışık, Piya Changmai, Nadin Rohland, Swapan Mallick, Nicole
549 Adamski, Deborah A Bolnick, Nasreen Broomandkhoshbacht, Francesca Candilio, Bren-
550 dan J Culleton, et al. Palaeo-eskimo genetic ancestry and the peopling of chukotka and
551 north america. *Nature*, 570(7760):236–240, 2019.
- 552 Matthew L Freedman, David Reich, Kathryn L Penney, Gavin J McDonald, Andre A
553 Mignault, Nick Patterson, Stacey B Gabriel, Eric J Topol, Jordan W Smoller, Carlos N
554 Pato, et al. Assessing the impact of population stratification on genetic association stud-
555 ies. *Nature genetics*, 36(4):388–393, 2004.
- 556 Xavier Grau-Bové, Sean Tomlinson, Andrias O O’Reilly, Nicholas J Harding, Alistair Miles,
557 Dominic Kwiatkowski, Martin J Donnelly, David Weetman, Anopheles gambiae 1000
558 Genomes Consortium, et al. Evolution of the insecticide target rdl in african anopheles
559 is driven by interspecific and interkaryotypic introgression. *bioRxiv*, pages 2019–12, 2020.
- 560 Christopher Heje Grønbech, Maximillian Fornitz Vording, Pascal N Timshel, Casper Kaae
561 Sønnerby, Tune Hannes Pers, and Ole Winther. scvae: Variational auto-encoders for
562 single-cell gene expression data. *bioRxiv*, page 318295, 2018.

REFERENCES

REFERENCES

- 563 Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David
564 Vazquez, and Aaron Courville. Pixelvae: A latent variable model for natural images.
565 *arXiv preprint arXiv:1611.05013*, 2016.
- 566 Ryan N Gutenkunst, Ryan D Hernandez, Scott H Williamson, and Carlos D Bustamante.
567 Inferring the joint demographic history of multiple populations from multidimensional
568 snp frequency data. *PLoS genet*, 5(10):e1000695, 2009.
- 569 Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew
570 Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual
571 concepts with a constrained variational framework. *Iclr*, 2(5):6, 2017.
- 572 Kent E Holsinger and Bruce S Weir. Genetics in geographically structured populations:
573 defining, estimating and interpreting f_{st} . *Nature Reviews Genetics*, 10(9):639–650, 2009.
- 574 Xianxu Hou, Linlin Shen, Ke Sun, and Guoping Qiu. Deep feature consistent variational
575 autoencoder, 2016.
- 576 Qiwen Hu and Casey S Greene. Parameter tuning is a key part of dimensionality reduction
577 via deep variational autoencoders for single cell rna transcriptomics. In *PSB*, pages 362–
578 373. World Scientific, 2019.
- 579 Thibaut Jombart. adegenet: a r package for the multivariate analysis of genetic markers.
580 *Bioinformatics*, 24(11):1403–1405, 2008.
- 581 Jerome Kelleher, Alison M Etheridge, and Gilean McVean. Efficient coalescent simu-
582 lation and genealogical analysis for large sample sizes. *PLoS Comput Biol*, 12(5):1–
583 22, 05 2016. doi: 10.1371/journal.pcbi.1004842. URL [http://dx.doi.org/10.1371%](http://dx.doi.org/10.1371%2Fjournal.pcbi.1004842)
584 [2Fjournal.pcbi.1004842](http://dx.doi.org/10.1371/journal.pcbi.1004842).
- 585 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*
586 *preprint arXiv:1412.6980*, 2014.
- 587 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
588 *arXiv:1312.6114*, 2013.
- 589 John Frank Charles Kingman. The coalescent. *Stochastic processes and their applications*,
590 13(3):235–248, 1982.
- 591 Dmitry Kobak and George C Linderman. Umap does not preserve global structure any
592 better than t-sne when using the same initialization. *bioRxiv*, 2019.
- 593 Tomasz Konopka. *umap: Uniform Manifold Approximation and Projection*, 2019. URL
594 <https://CRAN.R-project.org/package=umap>. R package version 0.2.3.1.
- 595 Maxime W Lafarge, Juan C Caicedo, Anne E Carpenter, Josien PW Pluim, Shantanu Singh,
596 and Mitko Veta. Capturing single-cell phenotypic variation via unsupervised representa-
597 tion learning. 2018.
- 598 Eric S Lander and Nicholas J Schork. Genetic dissection of complex traits. *Science*, 265
599 (5181):2037–2048, 1994.
- 600 Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther.
601 Autoencoding beyond pixels using a learned similarity metric, 2015.

REFERENCES

REFERENCES

- 602 Han Li and Peter Ralph. Local pca shows how the effect of population structure differs
603 along the genome. *Genetics*, 211(1):289–304, 2019.
- 604 Jianzhong Ma and Christopher I Amos. Investigation of inversion polymorphisms in the
605 human genome using principal components analysis. *PloS one*, 7(7), 2012.
- 606 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of*
607 *machine learning research*, 9(Nov):2579–2605, 2008.
- 608 Jonathan Marchini, Lon R Cardon, Michael S Phillips, and Peter Donnelly. The effects of
609 human population structure on large genetic association studies. *Nature genetics*, 36(5):
610 512–517, 2004.
- 611 Iain Mathieson and Gil McVean. Differential confounding of rare and common variants in
612 spatially structured populations. *Nature genetics*, 44(3):243, 2012.
- 613 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation
614 and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 615 Wes McKinney. Data structures for statistical computing in python. In Stéfan van der Walt
616 and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages
617 51 – 56, 2010.
- 618 Gil McVean. A genealogical interpretation of principal components analysis. *PLoS Genet*,
619 5(10):e1000686, 2009.
- 620 Paolo Menozzi, Alberto Piazza, and L Cavalli-Sforza. Synthetic maps of human gene fre-
621 quencies in europeans. *Science*, 201(4358):786–792, 1978.
- 622 Alistair Miles and Nick Harding. cggh/scikit-allele: v1.1.8, July 2017. URL <https://doi.org/10.5281/zenodo.822784>.
- 624 Alistair Miles, Nicholas J Harding, and the AG1000G consortium. Genetic diversity of the
625 African malaria vector *Anopheles gambiae*. *Nature*, 552(7683):96, 2017.
- 626 Daniel Mas Montserrat, Carlos Bustamante, and Alexander Ioannidis. Class-conditional
627 vae-gan for local-ancestry simulation. *arXiv preprint arXiv:1911.13220*, 2019.
- 628 John Novembre and Matthew Stephens. Interpreting principal component analyses of spatial
629 population genetic variation. *Nature genetics*, 40(5):646–649, 2008.
- 630 Travis Oliphant. NumPy: A guide to NumPy. USA: Trelgol Publishing, 2006–. URL
631 <http://www.numpy.org/>. [Online; accessed December 2019].
- 632 Nick Patterson, Alkes L Price, and David Reich. Population structure and eigenanalysis.
633 *PLoS genetics*, 2(12):e190, 2006.
- 634 Edzer Pebesma, Roger Bivand, Maintainer Edzer Pebesma, Suggests RColorBrewer, and
635 AAA Collate. Package ‘sp’. *The Comprehensive R Archive Network*, 2012.
- 636 Desislava Petkova, John Novembre, and Matthew Stephens. Visualizing spatial population
637 structure with estimated effective migration surfaces. *Nature genetics*, 48(1):94–100, 2016.

REFERENCES

REFERENCES

- 638 Cosimo Posth, Nathan Nakatsuka, Iosif Lazaridis, Pontus Skoglund, Swapan Mallick,
639 Thiseas C Lamnidis, Nadin Rohland, Kathrin Nägele, Nicole Adamski, Emilie Bertolini,
640 et al. Reconstructing the deep population history of central and south america. *Cell*, 175
641 (5):1185–1197, 2018.
- 642 Alkes L Price, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick,
643 and David Reich. Principal components analysis corrects for stratification in genome-wide
644 association studies. *Nature genetics*, 38(8):904, 2006.
- 645 Jonathan K Pritchard and Peter Donnelly. Case–control studies of association in structured
646 or admixed populations. *Theoretical population biology*, 60(3):227–237, 2001.
- 647 Michelle M Riehle, Tullu Bukhari, Awa Gneme, Wamdaogo M Guelbeogo, Boubacar
648 Coulibaly, Abdrahamane Fofana, Adrien Pain, Emmanuel Bischoff, Francois Renaud,
649 Abdoul H Beavogui, et al. The anopheles gambiae 2la chromosome inversion is associated
650 with susceptibility to plasmodium falciparum in africa. *Elife*, 6:e25813, 2017.
- 651 Noah A Rosenberg, Jonathan K Pritchard, James L Weber, Howard M Cann, Kenneth K
652 Kidd, Lev A Zhivotovsky, and Marcus W Feldman. Genetic structure of human popula-
653 tions. *science*, 298(5602):2381–2385, 2002.
- 654 Mashaal Sohail, Robert M Maier, Andrea Ganna, Alex Bloemendal, Alicia R Martin,
655 Michael C Turchin, Charleston WK Chiang, Joel Hirschhorn, Mark J Daly, Nick Pat-
656 terson, et al. Polygenic adaptation on height is overestimated due to uncorrected strati-
657 fication in genome-wide association studies. *Elife*, 8:e39702, 2019.
- 658 Dongfang Wang and Jin Gu. Vasc: dimension reduction and visualization of single-cell
659 rna-seq data by deep variational autoencoder. *Genomics, proteomics & bioinformatics*,
660 16(5):320–331, 2018.
- 661 Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec,
662 Francesco Montinaro, Cyril Furtlehner, Luca Pagani, and Flora Jay. Creating artifi-
663 cial human genomes using generative models. *bioRxiv*, 2019. doi: 10.1101/769091. URL
664 <https://www.biorxiv.org/content/early/2019/10/07/769091>.

665 Appendix 1

666 Other Things We Tried That Didn't Work

667 We tried a bunch of things while developing `popvae`. Here we document some of our dead-
668 ends in the hope they may be useful to others developing similar methods.

669 0.0.1 A Convolutional Neural Network

670 We first developed `popvae` using convolutional neural networks (CNNs) for both the encoder
671 and decoder. The feed-forward network we use here was originally intended as a naive
672 baseline for comparing our CNN performance, but it turned out to be faster and more
673 accurate (that is, lower validation loss), and had much lower memory requirements than
674 any CNN we tried. These included 2D CNNs run on phased haplotypes, 1D CNNs run on
675 unphased genotype counts, hybrid CNN+feed-forward networks stacking convolutional and
676 dense layers in succession, and restricting the CNN to either the encoder or the decoder.

677 0.0.2 A Recurrent Neural Network

678 We also tested recurrent neural networks (using the `cuDNNGRU()` layer in keras) as one or
679 both of the encoder/decoder pair. Due to memory limitations we were only able to test
680 relatively small, shallow networks with this approach (width 32, depth up to 3). Like the
681 CNNs these were slower, less accurate, and more resource-intensive than the dense network
682 we describe in the main text.

683 0.0.3 Skipping the Validation Set

684 It would be nice to not need a validation set. The train/test split introduces extra stochas-
685 ticity and you have to ignore some hard-earned data in training.

686 Unfortunately we couldn't find a good way of setting the learning rate scheduler or
687 establishing a good stopping time for model training without a validation set. Training
688 on all samples leads to constantly decreasing loss so all training runs go to the maximum
689 number of epochs. Examining the progress of latent spaces through model training for
690 these runs, the encoder seems to quickly identify and then refine structure in the input
691 samples, but eventually samples begin to cluster in a ring around the origin at 0,0. This
692 appears to reflect the Gaussian prior on the latent space dominating the loss function as the
693 reconstruction error approaches some lower bound. In runs with validation sets we observed
694 that validation loss typically increases once points begin circling the origin (Figure S1),
695 suggesting it reflects a typical overfitting behavior. Unfortunately the number of epochs
696 needed for this to occur is different for every dataset and we found no general solution for
697 estimating its location other than a validation set.

698 So we recommend using a validation set of at least 10%, and comparing latent spaces
699 from runs with multiple starting seeds (and so different train/validation splits). In a pinch,
700 users can set `--train_prop 1` to train on all samples and heuristically examine latent spaces
701 output during training to figure out a good stopping point.

702 0.0.4 Batch Normalization

703 Putting a batch normalization layer anywhere in either the encoder or decoder made vali-
704 dation loss worse in all of our tests.

REFERENCES

REFERENCES

705 **0.0.5 Dropout**

706 As above, dropout layers either made no difference or yielded slightly higher validation losses
707 no matter where we put it.

708 **0.0.6 Reweighting the Loss Function**

709 Higgins et al. (2017) proposed a modification of the standard VAE loss function which
710 amounts to multiplying the KL divergence by a factor β . This puts extra weight on the
711 normal prior of the latent space and on the MNIST dataset delivered more clustered and
712 interpretable latent spaces. Unfortunately the only suggested method for estimating β
713 in a truly unsupervised setting like ours is heuristic examination of model output. We
714 experimented with several values and found no consistent benefits either in latent space
715 or validation loss relative to our baseline approach. However this seems like a productive
716 area for further investigation and we plan to continue experimenting along these lines (and
717 encourage others to do so as well).

718 **Supplementary Figures and Tables**

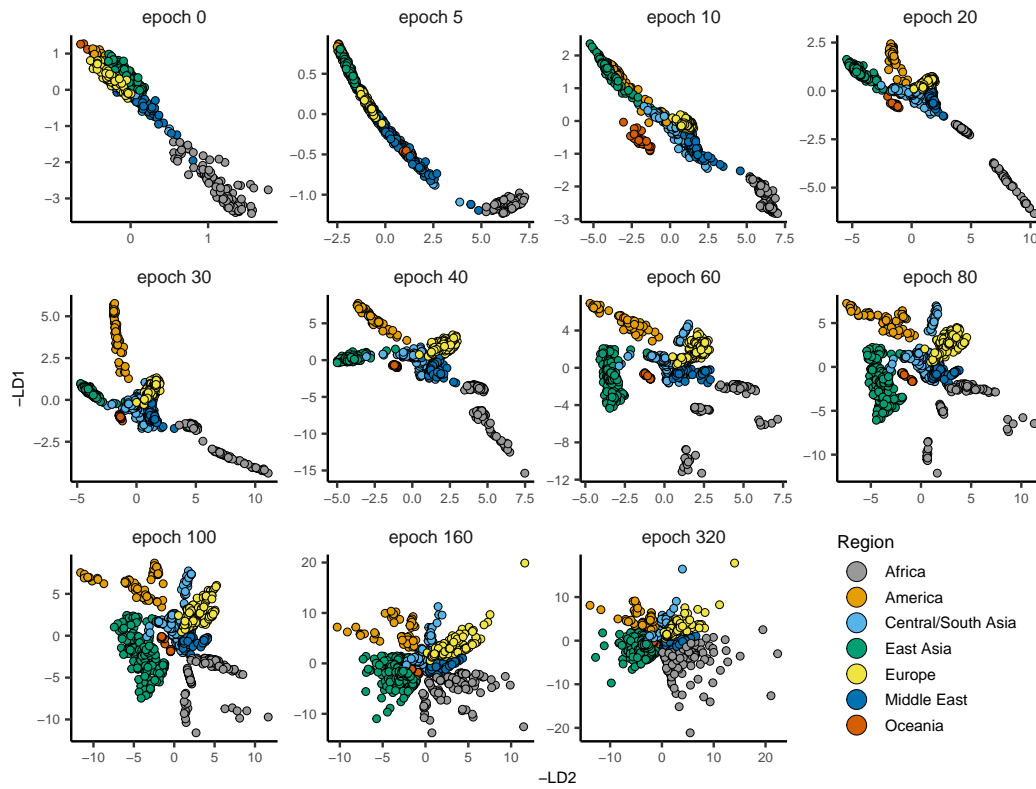


Figure S1: Latent spaces output during model training for HGDP data. Here `patience` was increased to 500 to show the overfitting behavior of `popvae`'s latent space. In this run validation loss was minimized at epoch 59.

REFERENCES

REFERENCES

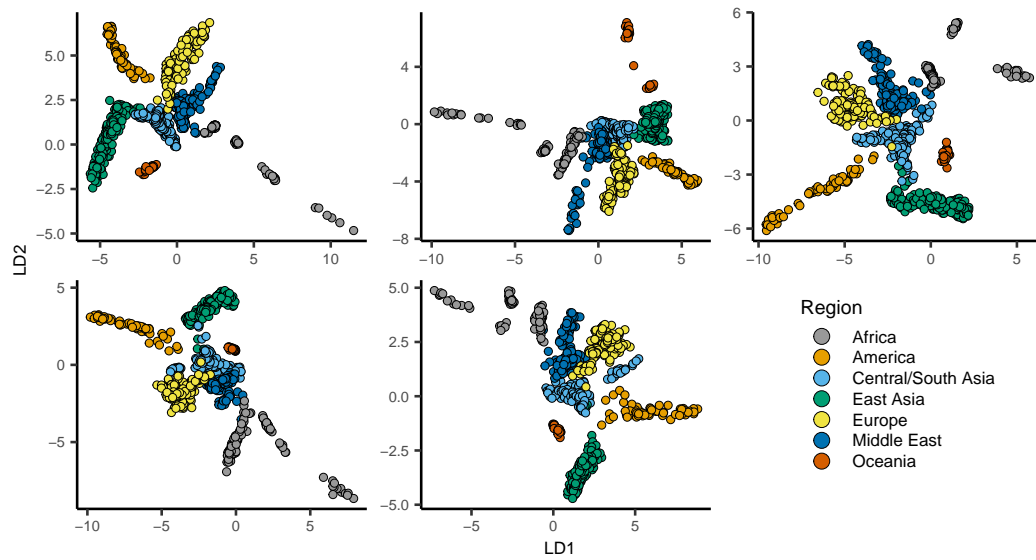


Figure S2: popvae latent spaces from runs with default hyperparameters and different random seeds.

REFERENCES

REFERENCES

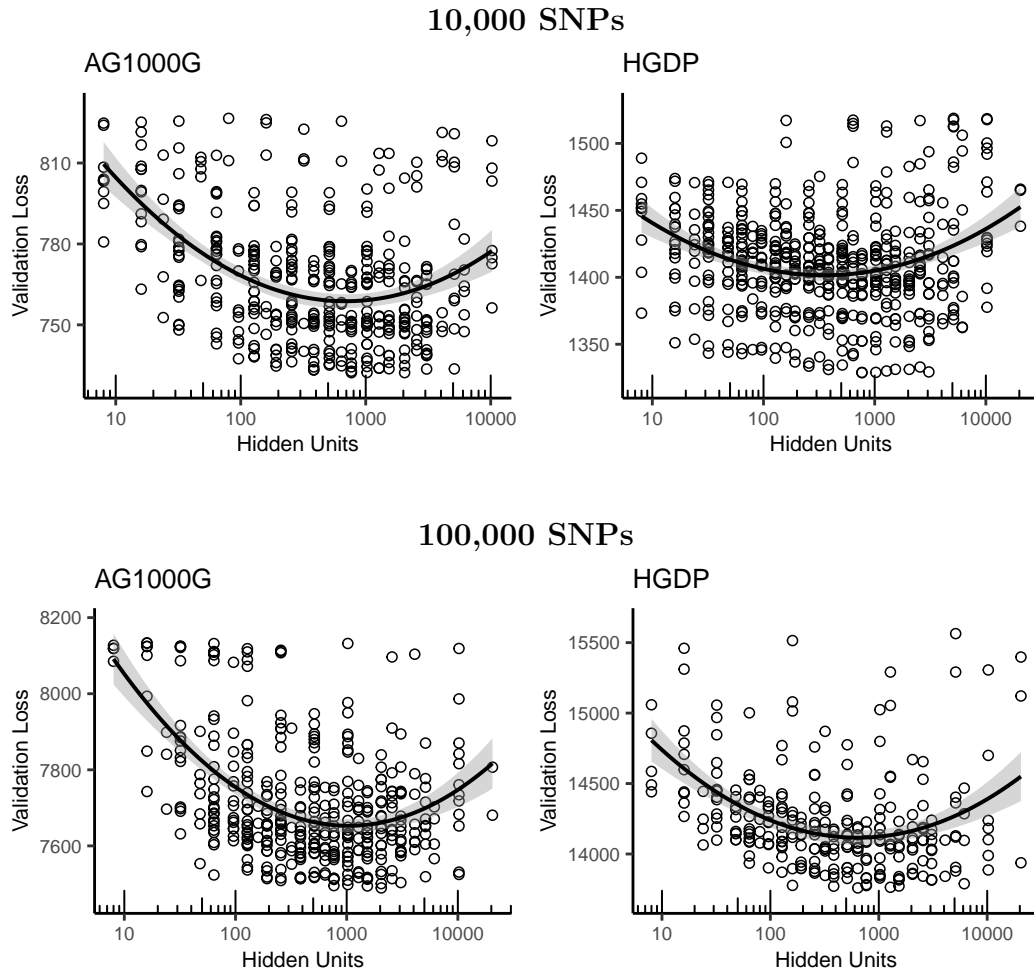


Figure S3: Validation loss as a function of the number of hidden units in a network for models fit to 100,000 SNPs selected randomly from *Anopheles* chromosome 3R and human chromosome 1. See table 1 for model rankings. Curves are quadratic least-squares model fits.

REFERENCES

REFERENCES

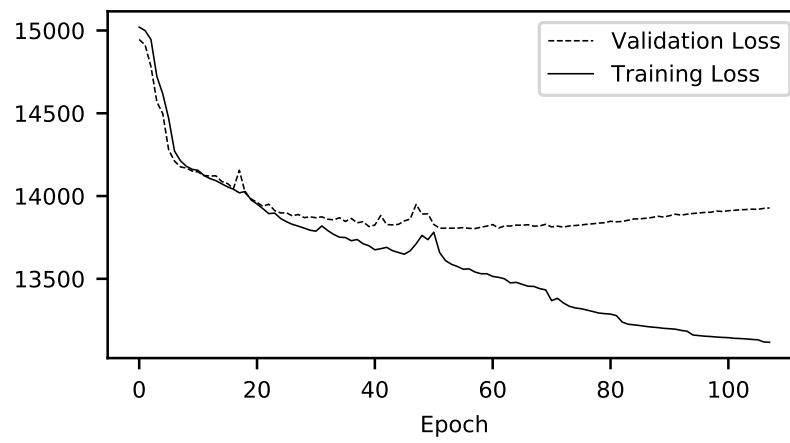


Figure S4: Example training history plot of showing training and validation loss by epoch during model training.

REFERENCES

REFERENCES

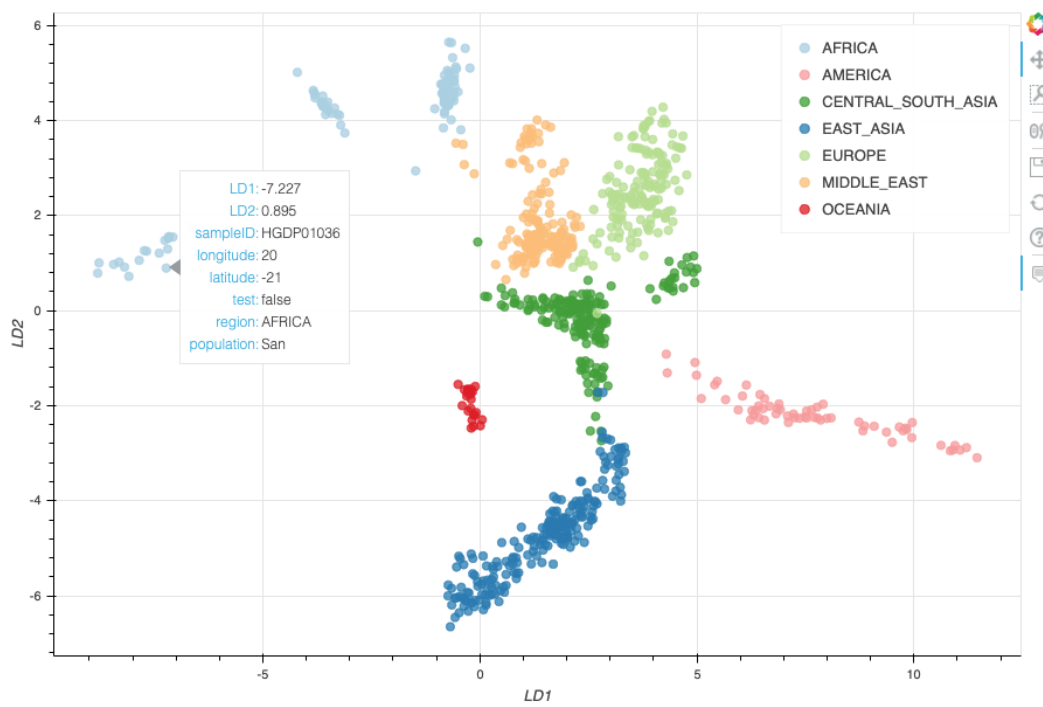


Figure S5: Example interactive plotting with scroll-over metadata.

REFERENCES

REFERENCES

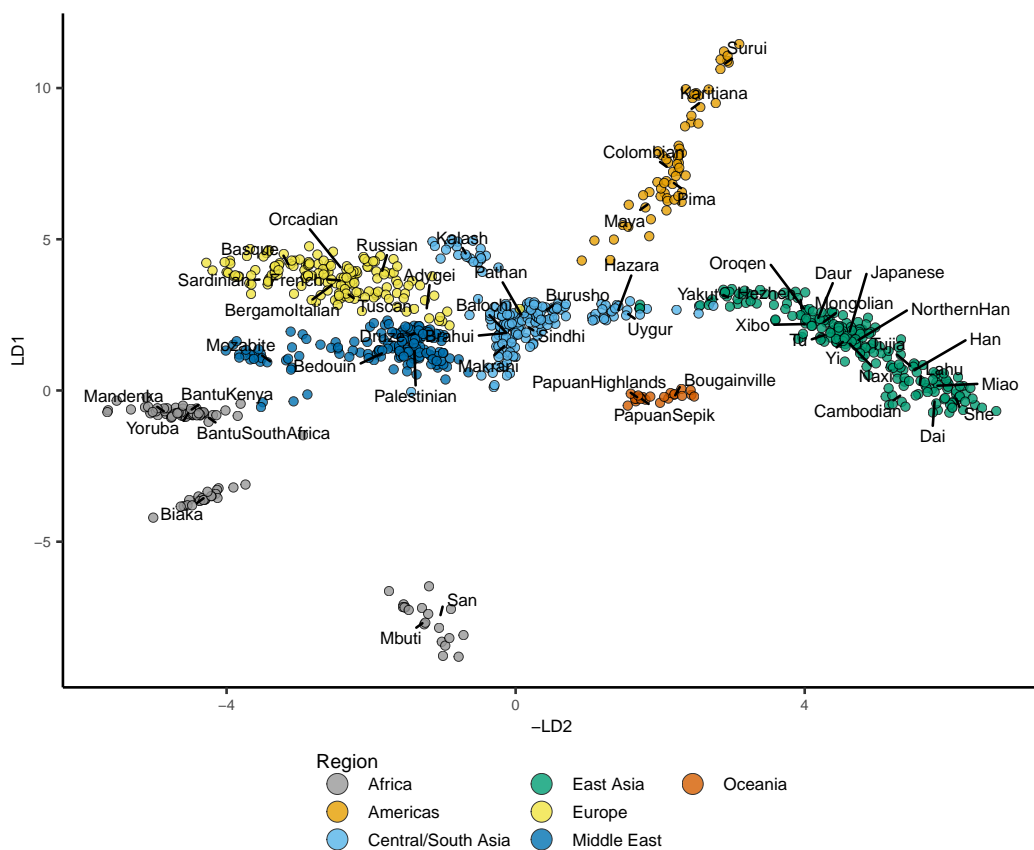


Figure S6: Latent space for 100,000 SNPs from chromosome 1 of the HGDP cohort (see Figure 2), with population centroids labeled.

REFERENCES

REFERENCES

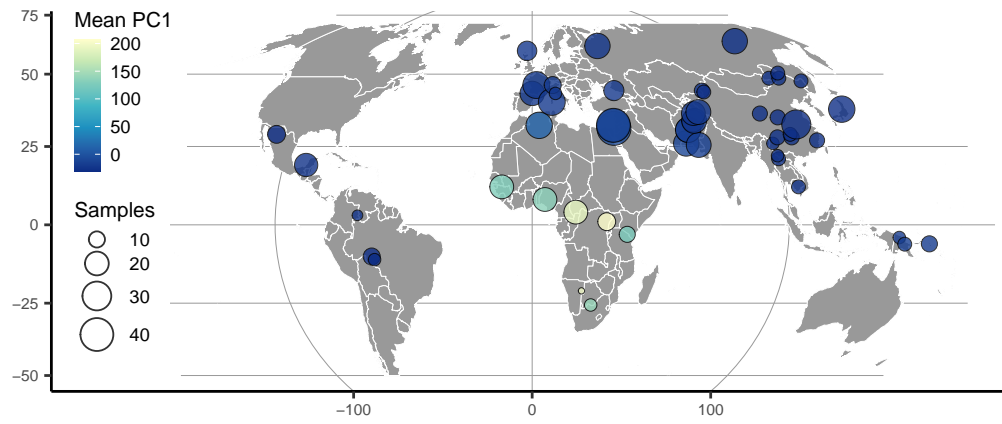


Figure S7: The first PC axis for HGDP SNPs summarized on a map as in Figure 3. Points show approximate population locations and are colored by the mean PC1 coordinate for each HGDP population. Densities show the distribution of PC1 scores for each HGDP region.

REFERENCES

REFERENCES

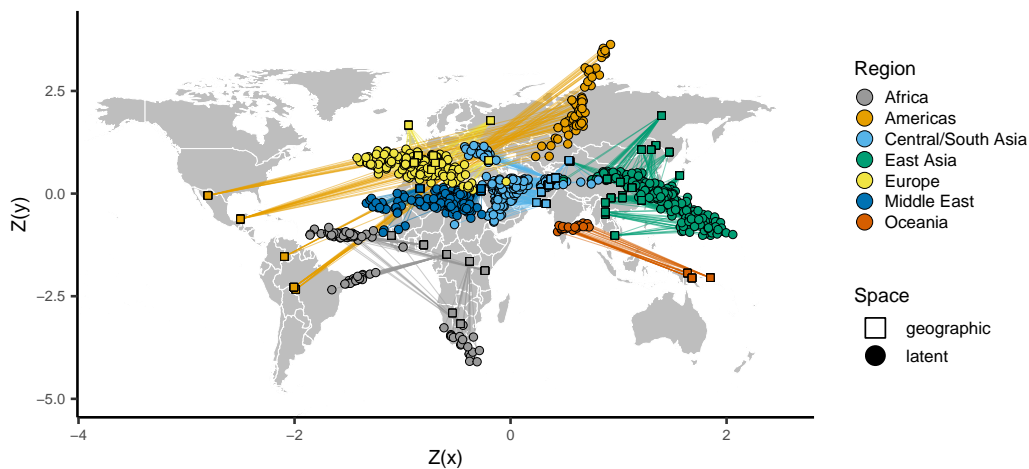


Figure S8: Comparing the VAE latent space with the geography of sampling localities HGDP samples. Circles show z-normalized sample locations in latent space and squares show the corresponding location in geographic space.

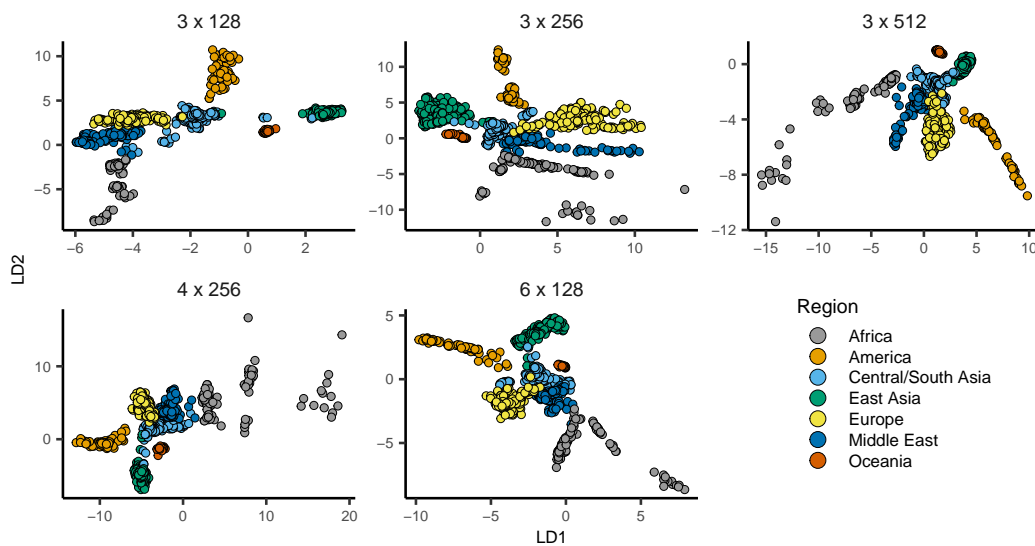


Figure S9: popvae latent spaces from runs with the same random seed and the top five network sizes by validation loss. Network sizes are listed as 'depth x width'.

REFERENCES

REFERENCES

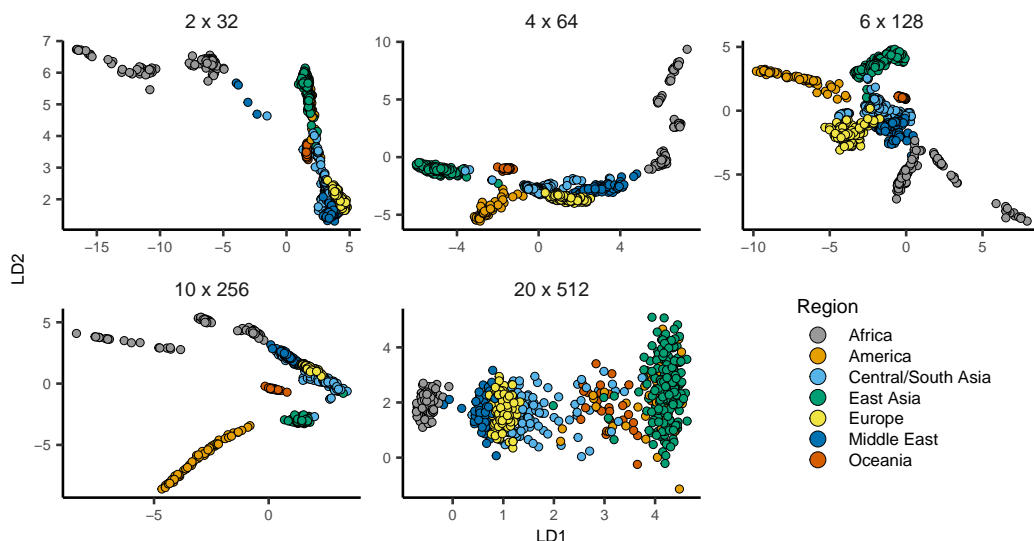


Figure S10: popvae latent spaces from models across the range of sizes tested. Network sizes are listed as ‘depth x width’.

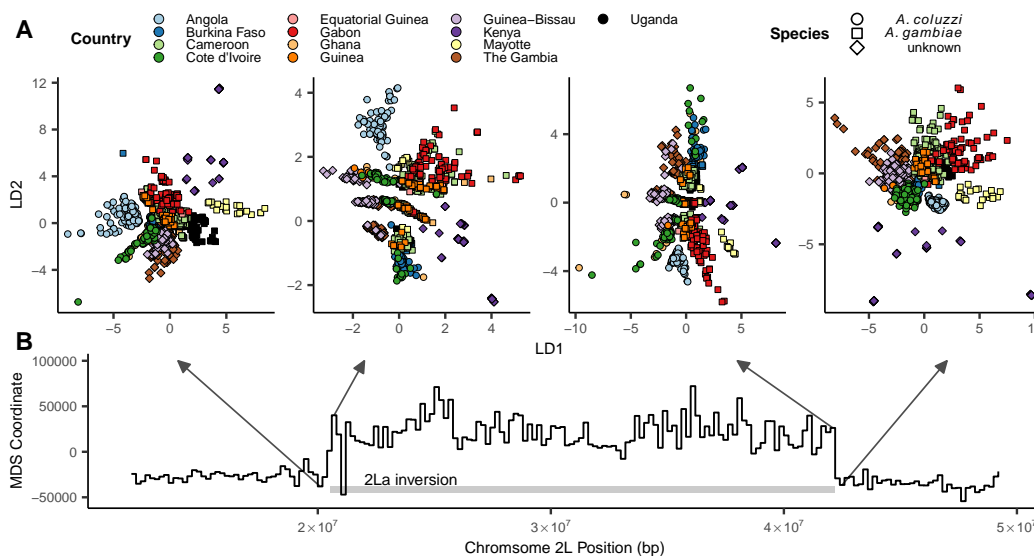


Figure S11: Latent spaces reflect inversion karyotypes at the 2La inversion in *A. gambiae* / *coluzzii*. A: VAE latent spaces for AG1000G phase 2 samples from windows near the 2La inversion breakpoints, with shapes indicating species and colors the country of origin. “Unknown” species localities include populations from Kenya, Guinea-Bissau, and the Gambia, for which diagnostic PCR markers are inconsistent or fail to amplify. B: Multi-dimensional scaling values showing difference in the relative position of individuals in latent space across windows – high values reflect windows in which samples cluster by inversion karyotype, and low values by species/region.

REFERENCES

REFERENCES

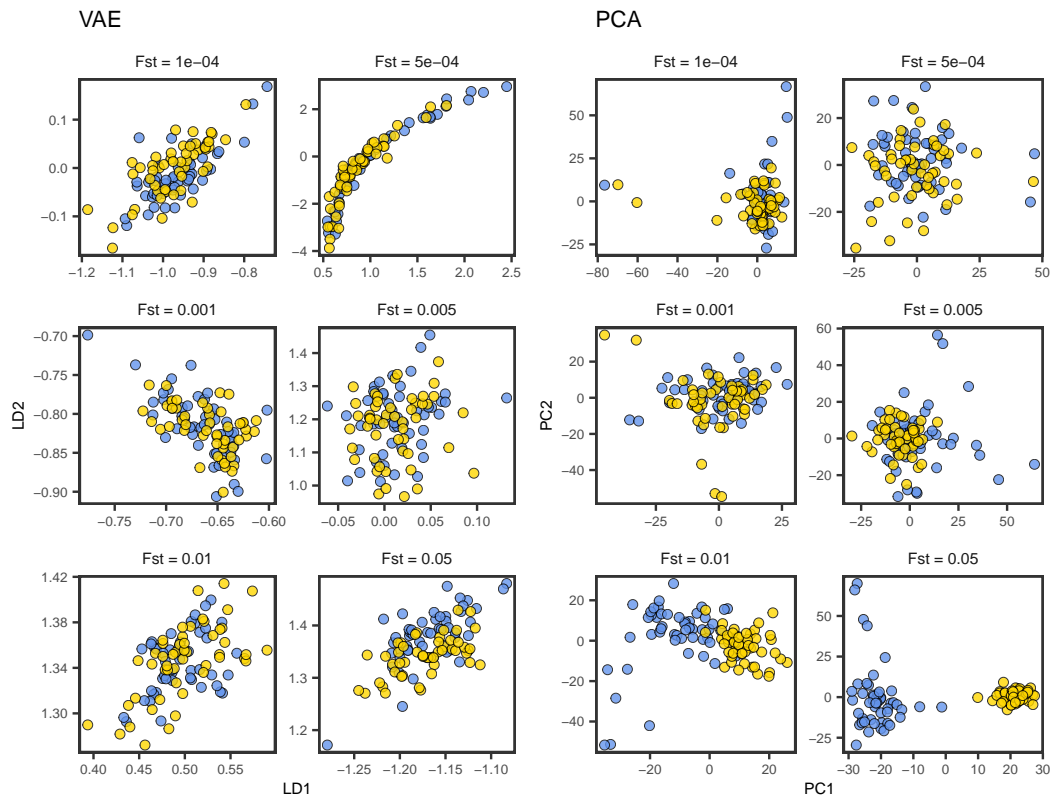


Figure S12: VAE latent spaces from the simulations shown in Figure 7, with `popvae` run at default settings.

REFERENCES

REFERENCES

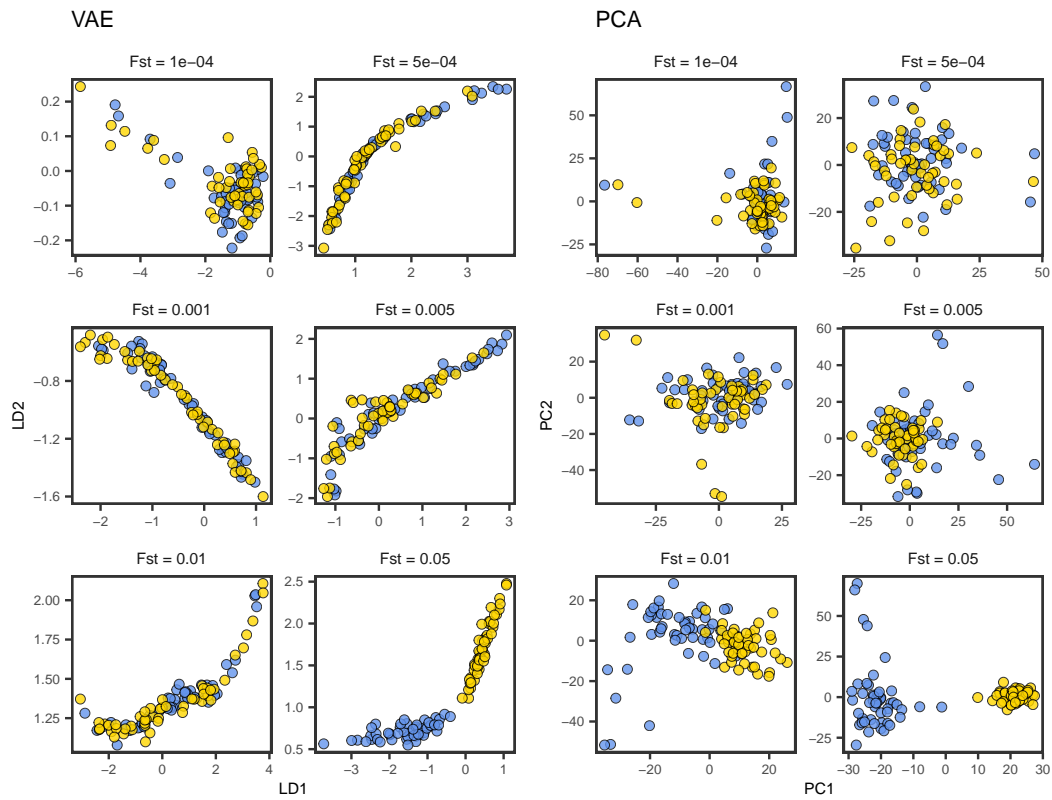


Figure S13: VAE latent spaces from the simulations shown in Figure 7, with `popvae` run with default network size (width 128, depth 6) and patience set to 500.

REFERENCES

REFERENCES

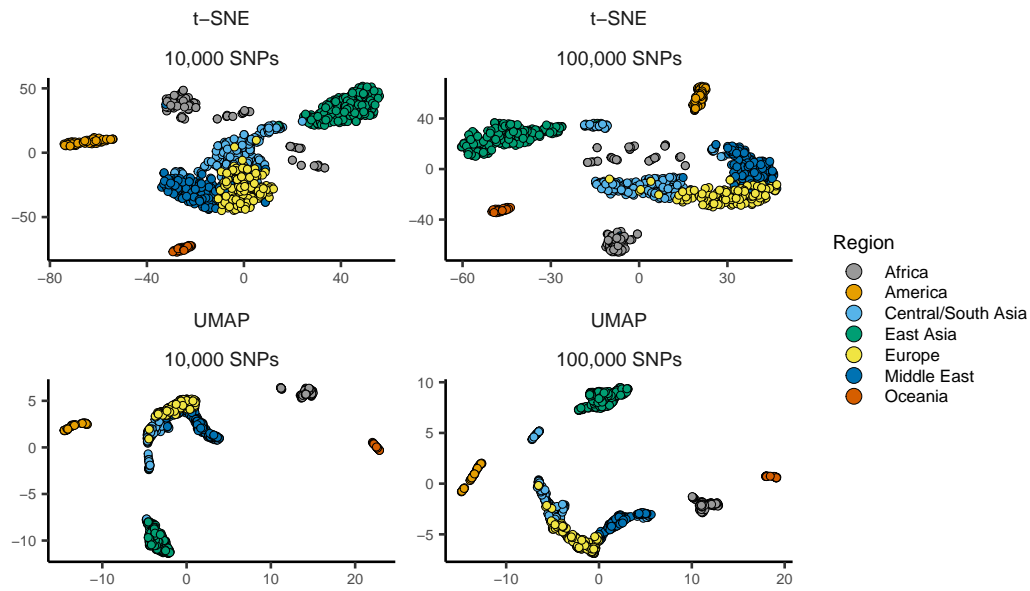


Figure S14: UMAP and t-SNE plots of HGDP samples using 100,000 or 10,000 SNPs. Both methods were run with default settings on the top 15 PC axes.

REFERENCES

REFERENCES

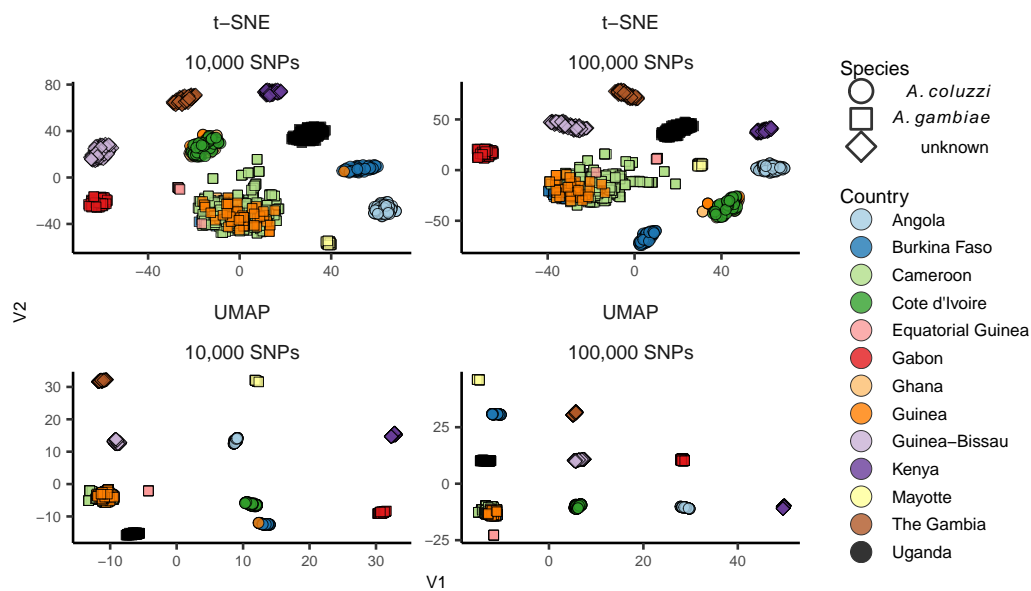


Figure S15: UMAP and t-SNE plots of AG1000G phase 2 samples using 100,000 or 10,000 SNPs at default settings. Both methods were run with default settings on the top 15 PC axes.

REFERENCES

REFERENCES

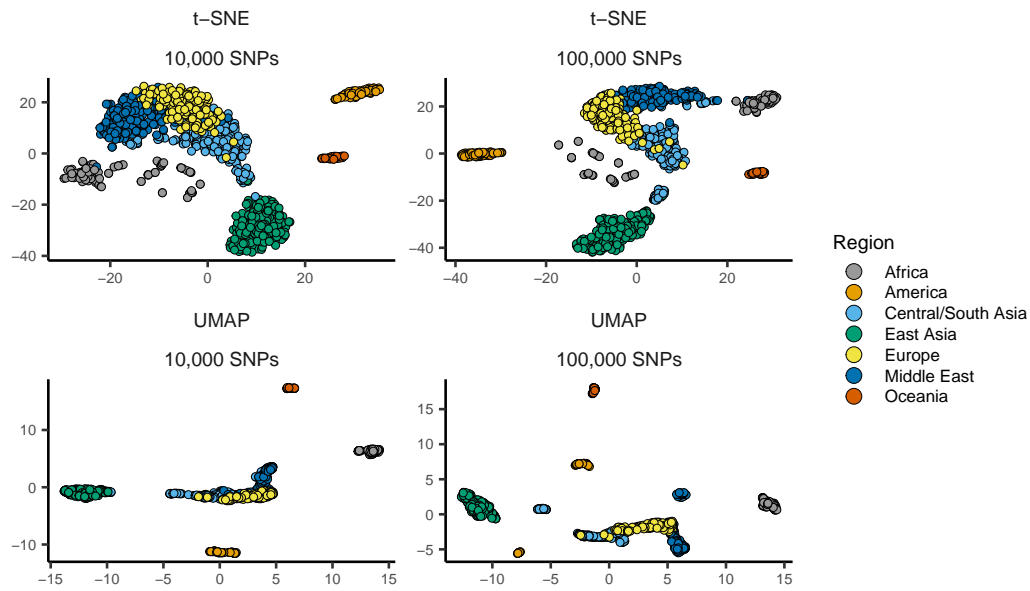


Figure S16: UMAP and t-SNE plots with parameters `n_neighbors=30` and `perplexity=60`. These settings are double the default values and are intended to improve global relative to local structure.

REFERENCES

REFERENCES

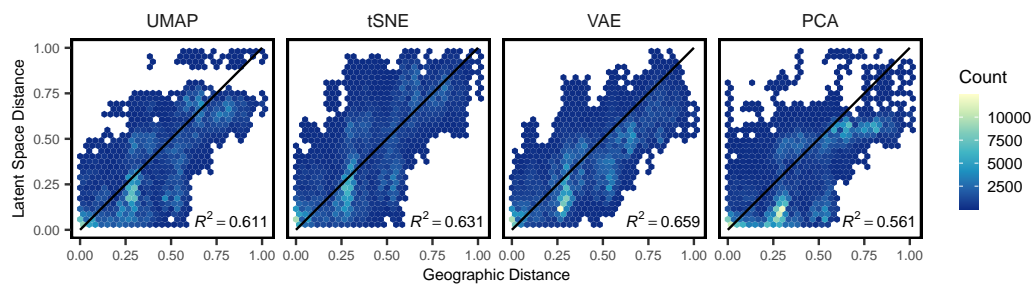


Figure S17: Comparison of relative pairwise distance for Eurasian HGDP samples, with UMAP parameter `n_neighbors=30` and t-SNE parameter `perplexity=60`. These settings are double the default values and are intended to improve global relative to local structure.

REFERENCES

REFERENCES

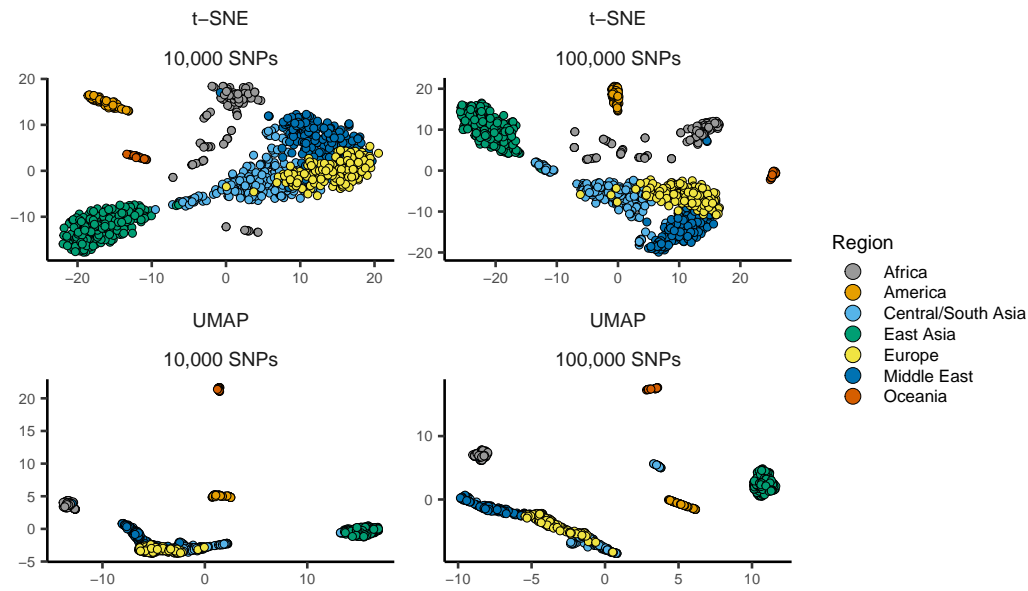


Figure S18: UMAP and t-SNE plots with parameters `n_neighbors=30` and `perplexity=60`. These settings are double the default values and are intended to improve global relative to local structure.

REFERENCES

REFERENCES

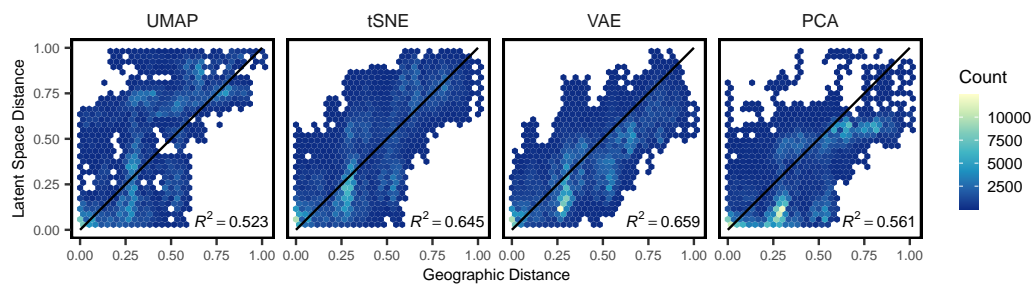


Figure S19: Comparison of relative pairwise distance for Eurasian HGDP samples, with UMAP parameter `n_neighbors=45` and t-SNE parameter `perplexity=90`. These settings are triple the default values and are intended to improve global relative to local structure.

REFERENCES

REFERENCES

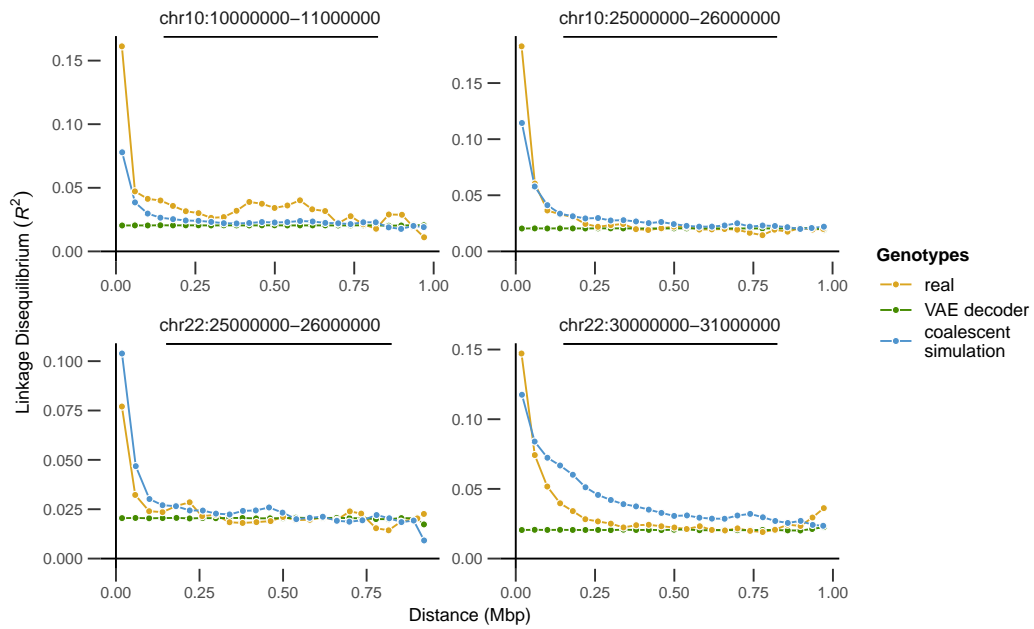


Figure S20: Comparing LD decay curves across real, simulated, and VAE decoder genotypes for four different regions of the genome. Points show the mean LD for all pairs of variants in each of 25 distance bins.

REFERENCES

REFERENCES

SNPs	HGDP			AG1000G		
	Depth	Width	Loss	Depth	Width	Loss
10,000	4	256	1394.231	3	256	750.677
	6	128	1394.48	6	128	750.859
	6	64	1394.504	4	128	751.0646
	10	64	1394.663	4	256	751.1514
	3	256	1394.976	6	64	751.7088
100,000	6	128	13955.76	4	256	7603.105
	3	256	13968.39	6	128	7606.528
	4	256	13971.75	3	256	7613.279
	3	512	13980.04	6	256	7614.232
	3	128	13992.27	4	128	7615.816
500,000	6	128	70087.32	10	128	37836.90
	10	64	70191.43	6	128	37848.67
	10	128	70203.22	6	64	37860.76
	6	64	70221.66	10	64	37872.49
	4	128	70357.73	4	128	37888.68

Table S1: Comparing validation loss across network sizes. Depth is the number of layers, width is the number of hidden units per layer, and loss is the mean validation loss across 5 random starting seeds for each network. Networks are ranked by loss for each dataset. SNPs were selected randomly from human chromosome 1 and *Anopheles* chromosome 3R.