



A deep learning framework for characterization of genotype data

Kristiina Ausmees* and Carl Nettelblad*,¹

*Department of Information Technology, Uppsala University, 752 37, Uppsala, Sweden

1

2 ABSTRACT

3 Dimensionality reduction is a data transformation technique widely used in various fields of genomics research.
4 The application of dimensionality reduction to genotype data is known to capture genetic similarity between
5 individuals, and is used for visualization of genetic variation, identification of population structure as well as
6 ancestry mapping. Among frequently used methods are PCA, which is a linear transform that often misses
7 more fine-scale structures, and neighbor-graph based methods which focus on local relationships rather than
8 large-scale patterns.
9 Deep learning models are a type of nonlinear machine learning method in which the features used in data
10 transformation are decided by the model in a data-driven manner, rather than by the researcher, and have
11 been shown to present a promising alternative to traditional statistical methods for various applications in
12 omics research. In this paper, we propose a deep learning model based on a convolutional autoencoder
13 architecture for dimensionality reduction of genotype data.

14 Using a highly diverse cohort of human samples, we demonstrate that the model can identify population clusters
15 and provide richer visual information in comparison to PCA, while preserving global geometry to a higher
16 extent than t-SNE and UMAP. We also discuss the use of the methodology for more general characterization
17 of genotype data, showing that models of a similar architecture can be used as a genetic clustering method,
18 comparing results to the ADMIXTURE software frequently used in population genetic studies.

19

KEYWORDS

deep learning, convolutional autoencoder, dimensionality reduction, genetic clustering, population genetics

1 INTRODUCTION

2 The increasing availability of large amounts of data has led to a rise
3 in the use of machine learning (ML) methods in several fields of
4 omics research. For many applications dealing with complex and
5 heterogeneous information, the data-driven approach has become
6 a promising alternative or complement to more traditional model-
7 based methods (Xu and Jackson 2019; Libbrecht and Noble 2015;
8 Schrider and Kern 2018).

9 Deep learning (DL) is an active subdiscipline of ML that has
10 had a large impact in several fields, including image analysis and
11 speech recognition (LeCun *et al.* 2015). DL methods comprise mod-
12 els that compute a nonlinear function of their input data using a
13 layered structure that learns abstract feature representations in a
14 hierarchical manner, and can be used for supervised learning tasks

such as prediction and also in unsupervised settings for pattern
recognition and data characterization problems (Goodfellow *et al.*
2016). A key aspect of DL is that the features used in data transfor-
mation are learned by the model as opposed to being defined by
the researcher, resulting in a higher level of flexibility than alter-
native ML algorithms such as support vector machines (Zou *et al.*
2019).

Advances have been made in developing DL techniques for
various types of omics data (Eraslan *et al.* 2019a). The current state-
of-the-art for predicting effects of genetic variants on splicing is
a DL model (Cheng *et al.* 2019). The DeepBind model (Alipanahi
et al. 2015) outperformed several previous non-DL approaches for
predicting sequence specificities of DNA-binding proteins. For the
task of variant calling of single-nucleotide polymorphisms (SNPs)
and small indels, the DeepVariant model of Poplin *et al.* (2018) was
shown to give improved results over existing tools (Nawy 2018).

DL has also been applied to unsupervised problems, in-

cluding imputation of metabolite and SNP data (Scholz *et al.* 2005; Chen and Shi 2019; Sun and Kardia 2008), de-noising of ChIP-sequencing data (Koh *et al.* 2017) and outlier detection of RNA sequencing gene expression data (Brechtmann *et al.* 2018). In the field of single-cell RNA-sequencing, DL methods have been used for imputation, de-noising as well as dimensionality reduction (Talwar *et al.* 2018; Ding *et al.* 2018; Eraslan *et al.* 2019b).

Dimensionality reduction is a data transformation technique that is commonly applied to SNP data in the fields of population and quantitative genetics. Applications include visualization of genetic variation, detection of population structure and correcting for stratification in genome-wide association studies (GWAS) (Patterson *et al.* 2006; Price *et al.* 2006). One of the most widely used methods for performing dimensionality reduction is Principal Component Analysis (PCA), in which a linear transformation is made onto uncorrelated dimensions that maximize the variance of the projected data (Pearson 1901). It has been shown that the lower-dimensional representation resulting from PCA can capture patterns in genetic variation, e.g. by reconstructing geographical relationships from genotype data (Novembre *et al.* 2008).

Although PCA is an efficient and reliable method, there are limitations associated with it. Firstly, it can be sensitive to attributes of sequence data such as the presence of rare alleles and SNPs that are correlated due to linkage disequilibrium (LD), which can cause groupings of samples that reflect such phenomena rather than genome-wide population structure (Ma and Shi 2020; Tian *et al.* 2008). To avoid such spurious effects, a stringent filtering procedure to remove low frequency variation and SNPs in high LD is usually required prior to performing PCA, although methods to handle LD by e.g. shrinkage methods have been proposed (Zou *et al.* 2010). Further limitations of PCA are related to the inability to capture nonlinear patterns in the data, as discussed in e.g. Alanis-Lobato *et al.* (2015), where the nonlinear method of non-centred Minimum Curvilinear Embedding (ncMCE) is proposed and shown to detect population structure in cases where PCA fails.

A family of nonlinear methods that has seen increased use on SNP data is neighbor graph-based models, including t-distributed stochastic neighbor embedding (t-SNE) and Uniform Manifold Approximation and Projection (UMAP) (van der Maaten and Hinton 2008; McInnes *et al.* 2020). These methods consider neighboring samples around each data point and try to find a lower-dimensional representation that preserves the distances between the points in the neighborhood. Both t-SNE and UMAP have been shown to be able to capture more fine-scale population structure than e.g. PCA, but the focus on preserving local topology results in a projection in which distances between larger clusters are more difficult to interpret (Gaspar and Breen 2019; Diaz-Papkovich *et al.* 2021).

More recently, DL methods for dimensionality reduction of SNP data have also been introduced. In Yelmen *et al.* (2021), the focus is on generating artificial genotypes, but one type of model considered, restricted Boltzmann machines (RBMs), projects the data to a reduced-dimensionality space which is compared to that of PCA. A DL approach based on variational autoencoders is presented in Battey *et al.* (2020), where they show that their model can capture subtle features of population structure, while preserving global geometry to a higher degree than both t-SNE and UMAP.

In this paper, we present a DL framework denoted Genotype Convolutional Autoencoder (GCAE) for nonlinear dimensionality reduction of SNP data based on convolutional autoencoders. The main differentiating feature between GCAE and the other DL meth-

ods mentioned is that our model makes use of convolutional layers, which take into account the sequential nature of genotype data. We describe adaptations to network architecture implemented to capture local as well as global patterns in sequence data, and compare dimensionality reduction performance to that of PCA, t-SNE and UMAP on a highly diverse cohort of human samples. We also demonstrate the broader applicability of the framework for general characterization of SNP data by showing that minor modifications in network structure can produce a model for solving the genetic clustering problem, and compare results to a model-based method commonly used in population genetic studies.

MATERIALS AND METHODS

Model architecture and training strategy

The proposed model for dimensionality reduction of SNP data is a convolutional autoencoder. Autoencoders are a class of DL models that transform data to a lower-dimensional latent representation from which it is subsequently reconstructed (Kramer 1991; Hinton and Salakhutdinov 2006). The idea is to learn features, or variables derived from the original data, that capture the important characteristics of the data. The structure of the model is shown in Figure 2. It comprises four types of layers that transform the input in a sequential manner: convolutional, pooling, fully-connected and upsampling layers.

Convolutional layers consist of a number of weight matrices, or kernels, of a specified size. These are used to compute a sliding dot product of the input. Each kernel is convolved over the input sequence along its spatial dimension with a given stride, or step size. In our models we use 1-dimensional convolutional layers with a stride of size 1. The depth of the layer output is thus determined by the number of kernels, and the spatial dimension of the data is unchanged. The pooling layers perform downsampling by applying a max filter over sliding windows of the data, separately for each depth dimension, reducing the size of the spatial dimension and leaving the depth unchanged.

The encoder alternates convolutional and pooling layers to increase the depth and reduce the spatial dimension of the data. The center of the model consists of a series of fully-connected layers in which the latent representation, or encoding, is defined. In contrast to convolutional layers, fully-connected layers contain weights between all pairs of variables in the input and output. The decoder roughly mirrors the structure of the encoder, with upsampling performed by means of nearest-neighbor interpolation to increase the spatial dimension of the data.

Residual connections, shown as black dashed lines in Figure 2, are used to stabilize the training process. These add the output from one layer to later parts of the network, skipping over layers in between, and have been shown to facilitate the optimization of deep networks for different applications (He *et al.* 2015).

Convolutional layers allow the model to capture local patterns and make use of the sequential nature of genetic data, allowing it to incorporate essential features such as LD at various length scales.

In order to facilitate the learning of global patterns in the input data, the model has two additional sets of variables. Each of these sets contains one variable per marker that is updated during the optimization process, allowing the model to capture marker-specific behavior. The two sets of marker-specific variables, illustrated in Figure 2 in red and green, are both inserted into the model by concatenation to layers in the decoder. One set of variables is also concatenated to the model input at every stage of the training process.

The activation function exponential linear unit (elu) is applied to convolutional layers, after which batch normalization is performed. The fully-connected layers also have elu applied to them, except the outermost ones which have linear activation. The final convolution is performed with a kernel size of 1, with linear activation and no batch-normalization. In order to regularize the network and avoid overfitting, dropout is used on the weights of the fully-connected layers, except those surrounding the latent representation, and Gaussian noise is added to the latent representation during training.

Input data is represented as an ($n_{samples} \times n_{markers}$) matrix of diploid genotypes, normalized to the range [0, 1] by mapping 0 → 0.0, 1 → 0.5, 2 → 1.0. When calculating the loss function, target genotypes are represented using one-hot encoding with 3 classes $[p(0), p(1), p(2)]$ with $p(g) = 1.0$ for the true genotype, and $p(g) = 0.0$ for the others. Model output o_{ij} of sample i at site j has the sigmoid function applied to it, after which it is interpreted as the allele frequency at the site. This scalar is transformed to the format of the targets using the formula for genotype frequencies $f(g)$ obtained from the principle of Hardy-Weinberg equilibrium: $[f(0) = (1 - o_{ij})^2, f(1) = 2(1 - o_{ij})o_{ij}, f(2) = o_{ij}^2]$.

The network is trained to reduce the categorical cross-entropy error (E) between target (y) and reconstructed (\hat{y}) genotypes, with an added L2 penalty on the values of the latent representation (e) for regularization. See Equation 1, where α is the regularization factor hyperparameter. Network optimization is performed by means of the ADAM algorithm (Kingma and Ba 2014), with a further exponential decay of the learning rate applied.

$$E(y, \hat{y}) = \sum_i^3 y_i \log(\hat{y}_i) + \alpha \sum_j^d e_j^2 \quad (1)$$

Additional regularization of the training process is performed by means of data augmentation. For every sample in the training process, a fraction of genotypes is randomly set to missing in the input data, represented by the value -1.0. A dimension representing missing and non-missing genotypes, with the values 0 and 1, respectively, is added to the input data, depicted in light blue in Figure 2. The fraction of missingness for each training batch is randomly selected from a pre-defined range. Noise was also added to the input data by introducing incorrect genotypes. With probability 0.2, missing genotypes of a batch were set to random genotype values, drawn from a uniform distribution.

The data was randomly split into training and validation sets consisting of 80% and 20% of samples, stratified by population. The training set was used for network optimization, with the termination criterion that the validation loss had not decreased for 300 epochs.

Different model architecture settings and hyperparameters were evaluated, and the best-performing setups were chosen by means of a hierarchical search procedure. We refer to Supplemental File S1 for details about the evaluated options, including the final model architecture settings and hyperparameter values used for obtaining the presented results.

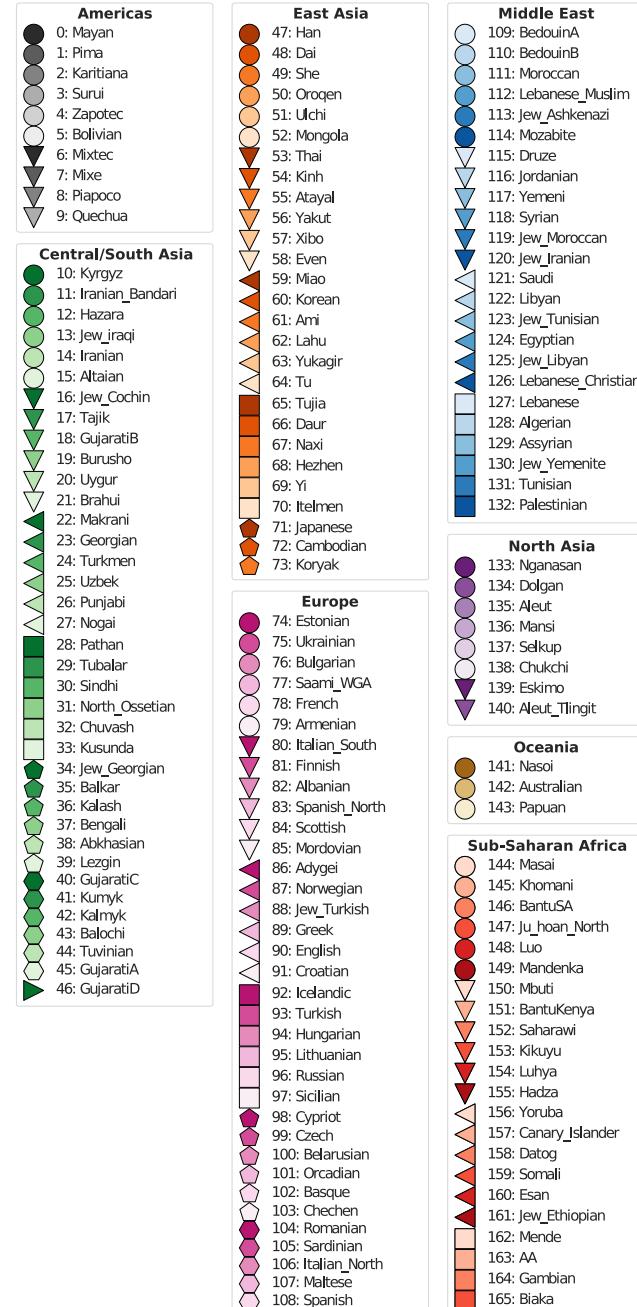


Figure 1 Populations and superpopulations of the Human Origins panel of genotype data. The coloring serves as a legend for Figure 3 and the numbering for Figures 5 and 6.

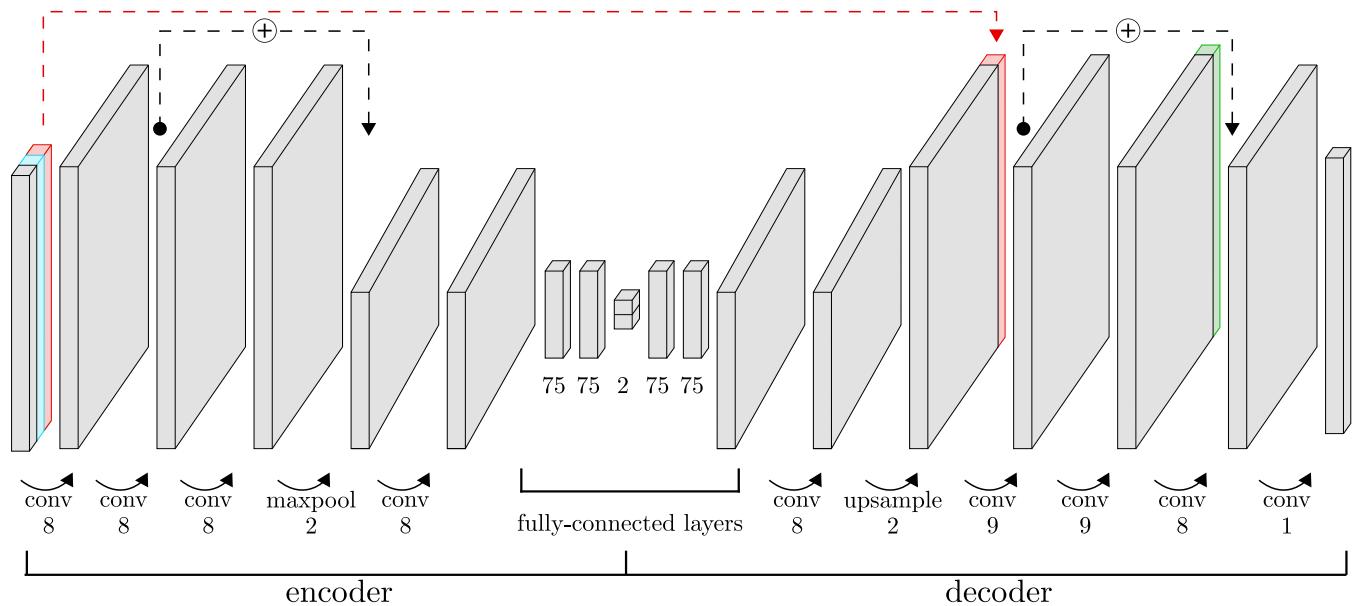


Figure 2 Architecture of the GCAE model used for dimensionality reduction. The encoder transforms data to a lower-dimensional latent representation through a series of convolutional, pooling and fully-connected layers. The decoder reconstructs the input genotypes. The input consists of three layers: genotype data (gray), a binary mask representing missing data (blue), and a marker-specific trainable variable per SNP (red). The red dashed line indicates where this marker-specific variable is concatenated to a layer in the decoder. Another marker-specific trainable variable, shown in green, is also concatenated to the second-last layer in the decoder. Black dashed lines indicate residual connections, where the output of a layer is added to that of another layer later in the network. The numbers below the layers indicate the number of kernels for convolutional layers, down- or upsampling factor for pooling and upsampling layers, and number of units for fully-connected layers. The displayed numbers are those of the final model used to obtain the presented results for dimensionality reduction to 2 dimensions. For other numbers of dimensions, the only modification made was to change the number of units in the latent representation from 2 to 4, 6, 8 or 10. For the genetic clustering application, the number of units in the latent representation was $k = 5$.

1 Human Origins data set

2 The data set used to evaluate dimensionality reduction and ge-
3 netic clustering performance was derived from the fully public
4 present-day individuals of the Affymetrix Human Origins SNP
5 array analyzed in Lazaridis *et al.* (2016). This data set is designed
6 for population genetic studies and represents worldwide genetic
7 variation, containing 2,068 samples from 166 populations. These
8 were categorized into 8 superpopulations, displayed in Figure 1
9 which also serves as a legend for Figures 3, 5 and 6.

10 The data was filtered to exclude sex chromosomes and non-
11 informative sites, and one sample (NA13619) was removed due to
12 relation to another (HGDP01382). In order to obtain a single
13 data set for fair comparison between methods, the genotypes were
14 further filtered according to the procedure that is common to per-
15 form prior to applying PCA on SNP data. A minor allele frequency
16 (MAF) threshold of 1% was enforced, and LD pruning was per-
17 formed by removing one of each pair of SNPs in windows of 1
18 centimorgan that had an allelic R^2 value greater than 0.2.

19 As the comparison of robustness of different methods to missing
20 data was beyond the scope of this study, missing genotypes were
21 set to the most frequent value per SNP so as to avoid their influence
22 over dimensionality reduction results. The final data set consisted
23 of 2,067 individuals typed at 160,858 biallelic sites.

24 Evaluation of dimensionality reduction performance

25 Comparison of performance between GCAE, PCA, t-SNE and
26 UMAP was performed by means of evaluating the ability of the
27 dimensionality reduction to capture population structure. A k -
28 Nearest Neighbors (k -NN) classification model was defined based
29 on the projected data by assigning a population label to each sam-
30 ple based on the most frequent label among its k nearest neighbors.

31 The evaluation was performed using 2, 4, 6, 8, and 10 latent
32 dimensions. For UMAP, t-SNE and GCAE, hyperparameter tun-
33 ing was done for each number of latent dimensions, selecting the
34 configuration that yielded the highest F1 score for a 3-NN classi-
35 fication model. See Supplemental File S1 for details. The model
36 with the selected hyperparameters was then used to calculate F1
37 scores for classification models using 3 and 20 neighbors. The
38 different numbers of neighbors were used to obtain metrics that
39 capture different aspects of performance, e.g. tightness and degree
40 of overlap of the populations clusters.

41 Further, in order to evaluate the ability of the models to capture
42 global patterns in the data, the models were also evaluated on their
43 ability to classify samples according to membership to the larger
44 continental groups. This was done by calculating the F1 scores
45 using the superpopulations, rather than populations, as labels (see
46 Figure 1).

47 Classification performance was measured by the F1 score, de-
48 fined per class label c as the harmonic mean of the precision and
49 recall: $F1_c = \frac{2 * precision * recall}{precision + recall}$. The total F1 score for a model was de-
50 fined as the average F1 score, weighted by the number of samples
51 per class.

52 PCA was performed using the same normalization method as in
53 the software SMARTPCA (Patterson *et al.* 2006), i.e. by subtraction
54 of the mean and division with an estimate of the standard deviation
55 of the population allele frequency per SNP. For t-SNE and UMAP,
56 the data was standardized by removing the mean and scaling to
57 unit variance.

58 The entire data set was used for performance evaluation. This
59 is standard practice for using PCA, t-SNE and UMAP, particularly
60 the latter two as they are non-parametric models. For neural net-
61 works, a test set of previously unseen samples is usually used for

62 performance evaluation. In this case, we motivate the use of the
63 entire data set by the nature of the research question itself and
64 the fact that the performance metrics are based on the population
65 labels which are unseen by the network and therefore not used in
66 the model optimization.

67 Extension to genetic clustering

68 The problem of genetic clustering refers to the characterization of
69 individual genomes by proportional assignment to a set of clusters,
70 or genetic components. These may be used in the analysis of pop-
71 ulation structure and in identifying patterns of genetic variation
72 between populations.

73 The DL model for genetic clustering was developed by making
74 minor changes to the autoencoder architecture used for dimen-
75 sionality reduction described above. The number of units in the
76 encoding layer was changed from 2 to k , the number of clusters. In
77 order to obtain a proportional assignment, softmax normalization
78 was further applied to the encoding to obtain a vector of k values
79 that sum to 1. The loss function was the same as for the dimension-
80 ality reduction application, shown in Equation 1. Supplemental
81 File S1 contains more information about the network architecture
82 settings and training options used for the genetic clustering model.

83 We consider the widely used software ADMIXTURE (Alexander-
84 der *et al.* 2009) as a comparison method for the genetic clustering
85 application, and present results in a similar manner using bar
86 graphs displaying the proportional assignment of clusters for each
87 sample. For both GCAE and ADMIXTURE, the Human Origins
88 data set described above was used.

89 Analysis of output genotypes

90 In order to further analyze the representation learned by the model,
91 we compared different characteristics of the output genotypes to
92 those of the true ones. First, we studied the distributions of allele
93 frequencies by comparison of the respective site frequency spectra.
94 Secondly, we compared the spatial structure in the two data sets
95 by studying the pattern of LD decay along the chromosome.

96 For this analysis, we considered a different data set consisting
97 of a single chromosome typed at a denser set of SNPs. As in Battey
98 *et al.* (2020), we used chromosome 22 from the 1000 Genomes phase
99 3 data (Auton *et al.* 2015), but restricted to biallelic SNPs in the
100 region 24500000-26500000 bp, resulting in 61104 sites.

101 Network optimization was performed using all 2504 samples,
102 randomly split into training and validation sets consisting of 80%
103 and 20% of samples, stratified by population. The same training
104 procedure as described in was used, as well as a similar 2-D model
105 architecture, with the difference that a larger kernel size was used
106 for the convolutional layers. See Supplemental File S1 for details
107 on the model and hyperparameters evaluated. Model evaluation
108 was done based on the one that yielded the most non-fixed sites
109 for use in the LD calculation. For these experiments, we also used
110 weighting of the loss function to handle the skewed distribution of
111 genotypes, by means of a class-balanced loss based on the effective
112 number of samples (Cui *et al.* 2019) with $\beta = 0.95$.

113 LD analysis was performed on the output genotypes of the
114 trained GCAE model, on a subset of samples and SNPs. Similar
115 to the LD analysis performed in Battey *et al.* (2020), we considered
116 samples from the YRI population only, and SNPs in the interval
117 25000000-26000000 bp. We further applied a MAF threshold of 1%
118 and restricted the genotypes to only include those passing all filters
119 in the 'strict' accessibility mask provided by the 1000 Genomes
120 Project.

121 The R^2 measure of LD was calculated for the output genotypes

and compared to that of the true genotypes. As R^2 is not defined for pairs of loci where one or more allele frequencies are equal to zero, sites for which only one allele was present in the output genotypes were excluded from the analysis, in order to consider the same set of sites in the true and output genotypes.

Implementation

GCAE is implemented in Python 3 using Tensorflow 2 (Abadi *et al.* 2015) and is available at <https://github.com/kausmees/GenoCAE> as a command-line program. The programs plink 1.9 (Shaun Purcell, Christopher Chang 2020; Chang *et al.* 2015) and bcftools 1.14 (Danecek *et al.* 2021) were used for filtering of genotype data. All other preprocessing of data, as well as performance evaluation and visualization, was implemented in Python. The library scikit-learn 1.0 (Pedregosa *et al.* 2011) was used for calculating the F1 score metric, with scikit-allel 1.3.5 (Miles *et al.* 2021) being used for LD analysis.

PCA and t-SNE were performed using the Python libraries scikit-learn and MulticoreTSNE 0.1 (Ulyanov 2016), respectively. The reference results for genetic clustering were obtained using the software ADMIXTURE 1.3.0 using the em method.

CPU computations were performed on the resources of Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) on a cluster of compute servers equipped with 128 GB memory, each comprising two 8-core Xeon E5-2660 processors. GPU computations were run on National Supercomputer Centre (NSC) at Linköping University, on an NVIDIA SuperPOD with DGX-A100 nodes equipped with 8 NVIDIA A100 Tensor Core GPUs with 40 GB on-board HBM2 VRAM, 2 AMD Epyc 7742 CPUs, 1 TB RAM.

RESULTS

Figure 3 shows dimensionality reduction results using GCAE, PCA, t-SNE and UMAP on the Human Origins data set. On a global scale, PCA and GCAE result in similar patterns. Both methods are able to capture global geometry to a high degree, with a visible clustering according to superpopulation. Both methods also result in a consistent global pattern with the Sub-Saharan African superpopulation separating distinctly, a gradient from the Middle East and Europe to South/Central and North Asia onto East Asia that roughly reflects the geography of Eurasia, and a separate Oceanian cluster.

A difference is that while PCA essentially clusters all non-African populations on one dimension, the GCAE plot is more spread out, with a north-south gradient in addition to the east-west relationships mainly captured by PCA. This is visible within superpopulations as well as between them, e.g. by the appearance of the Americas as a distinct cluster and a clearer differentiation of North Asia from and Central/South and East Asia. Populations of the Far East Siberia like Yukagir, Koryak, Chukchi and Eskimo, for example, appear clearly set apart along the D1 axis, which is also the dimension that mainly distinguishes the American samples. Within the African superpopulation, samples from the north and east of the continent are in closer proximity to the Middle Eastern cluster, and differentiated from the cluster of west-African populations. The San populations Khomani and Ju Hoan also separate distinctly, which is not evident in the PCA plot. A similar observation holds for the Mbuti and Biaka groups of the Congo Basin area.

The neighbor graph-based methods, in contrast, do not display the same global pattern as PCA and GCAE. While t-SNE also shows clustering according to superpopulation, it is a different

pattern with Sub-Saharan Africa in the middle, South/Central Asian populations more spread out, and some Asian populations in tight formations at large distances from the main cluster. UMAP shows even more populations appearing in tight and highly separated clusters, with the rest forming a distinct shape on a curved line. While more difficult to interpret, some degree of clustering according to superpopulation is also visible. Overall, both t-SNE and UMAP result in visualizations with less correspondence to global geographical patterns, but a comparatively high degree of clustering of individual populations.

It is important to note that PCA used as a dimensionality reduction method differs from the other methods considered in that there is a choice of which principal components to use. In this part of the evaluation, the focus is on the visual information in the two-dimensional projection, for which it is common to use the first two principal components only. We refrain from discussing other combinations of principal components here, but note that there may be different structure visible when selecting others. For a more complete assessment, we refer to the results of the classification performance below, which indicate the ability of the models to capture population structure when considering multiple dimensions.

Figure 4 shows the F1 scores of 3-NN and 20-NN classification models based on the dimensionality reduction of GCAE, PCA, TSNE and UMAP for 2-10 dimensions. The top and bottom plots show scores for the population and superpopulation classification models, respectively.

For the population classification model, UMAP resulted in the highest F1 scores, which is consistent with the tight clustering of individual populations in Figure 3. t-SNE tends to give relatively high classification performance for lower dimensions, and does not show much improvement in score with increased dimensionality, a trend that is also visible for UMAP. PCA and GCAE show a similar pattern of increasing performance, with GCAE tending to have higher scores.

For the superpopulation classification model, the relative performance of PCA and GCAE to the other models increases. This indicates that the neighbor graph-based methods have less ability to capture global structures, and is particularly evident in the 20-NN model for which they have significantly lower scores. PCA shows quite consistent performance gains for increased dimensionality, indicating that the additional PCs do add structure that is useful for the classification model. For lower numbers of dimensions, GCAE tends to have higher scores than PCA, with performance dropping off for higher dimensions. A possible explanation for this is that regularizing the latent space of GCAE becomes more difficult with increased size, and further exploration of hyperparameter space and regularization methods might be needed to utilize the space to a larger extent.

It is worth noting that our optimization criteria for hyperparameters clearly favors local performance. Both t-SNE and UMAP are models for which hyperparameters control the balance between attention given to local and global aspects of the data, and the authors of UMAP argue that their model can achieve a higher degree of preservation of global patterns in comparison to t-SNE (McInnes *et al.* 2020). The hyperparameter selection process here does not necessarily reflect that of a user interested in using t-SNE or UMAP for visualization purposes, where a more subjective evaluation would most likely be used. To give a more comprehensive view, we evaluate the visualization and clustering results of other hyperparameter combinations for UMAP in Supplemental File S1.

Figures 5 and 6 show the genetic clustering results on the Hu-

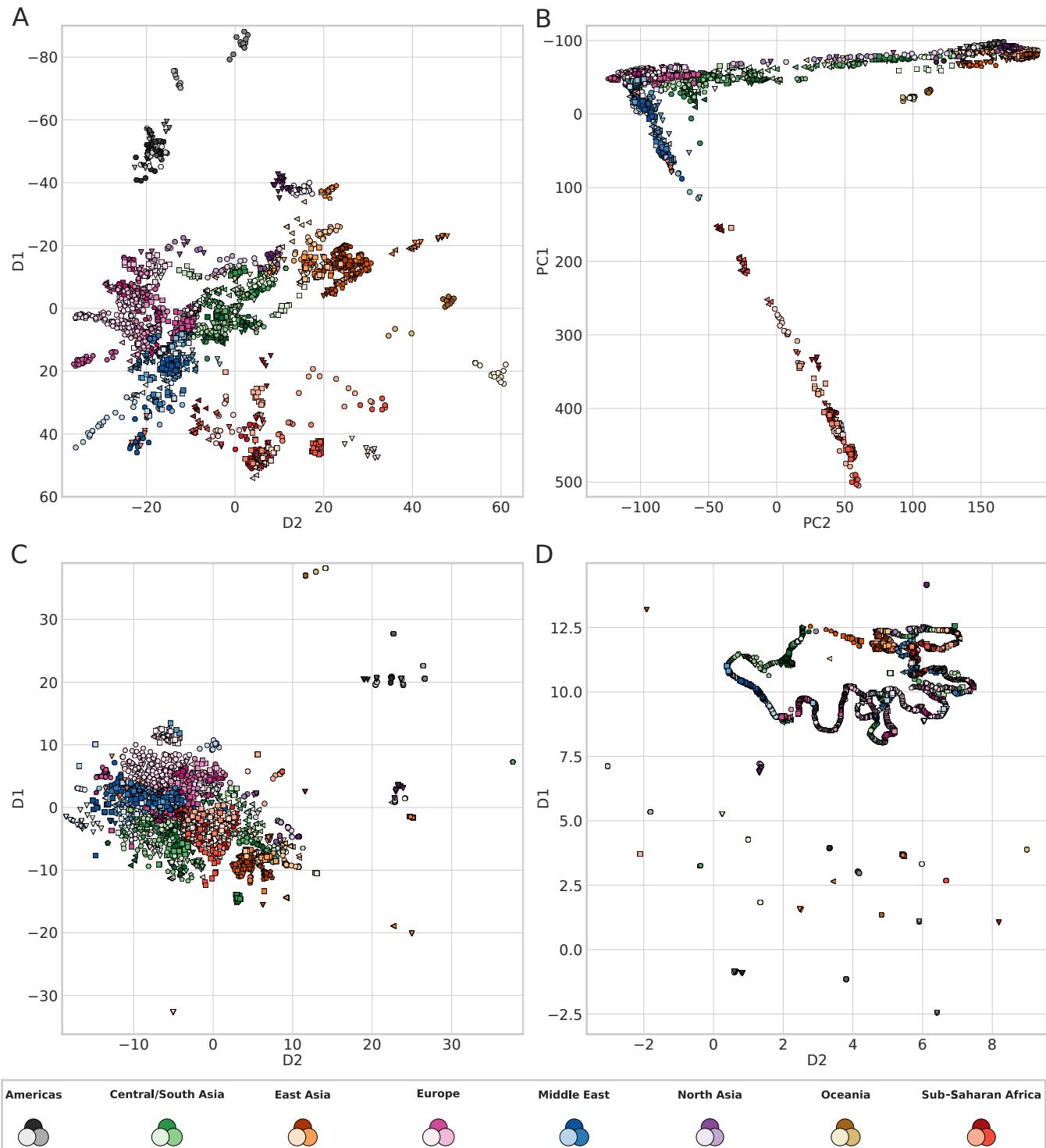


Figure 3 Dimensionality reduction results for GCAE (A), PCA (B), t-SNE (C) and UMAP (D) on the Human Origins data set. For GCAE and PCA, the D1 and PC1 axes have been inverted in order to get a more direct correspondence to the cardinal geographical directions. Legend shows superpopulation colors, full legend with all populations in Figure 1.

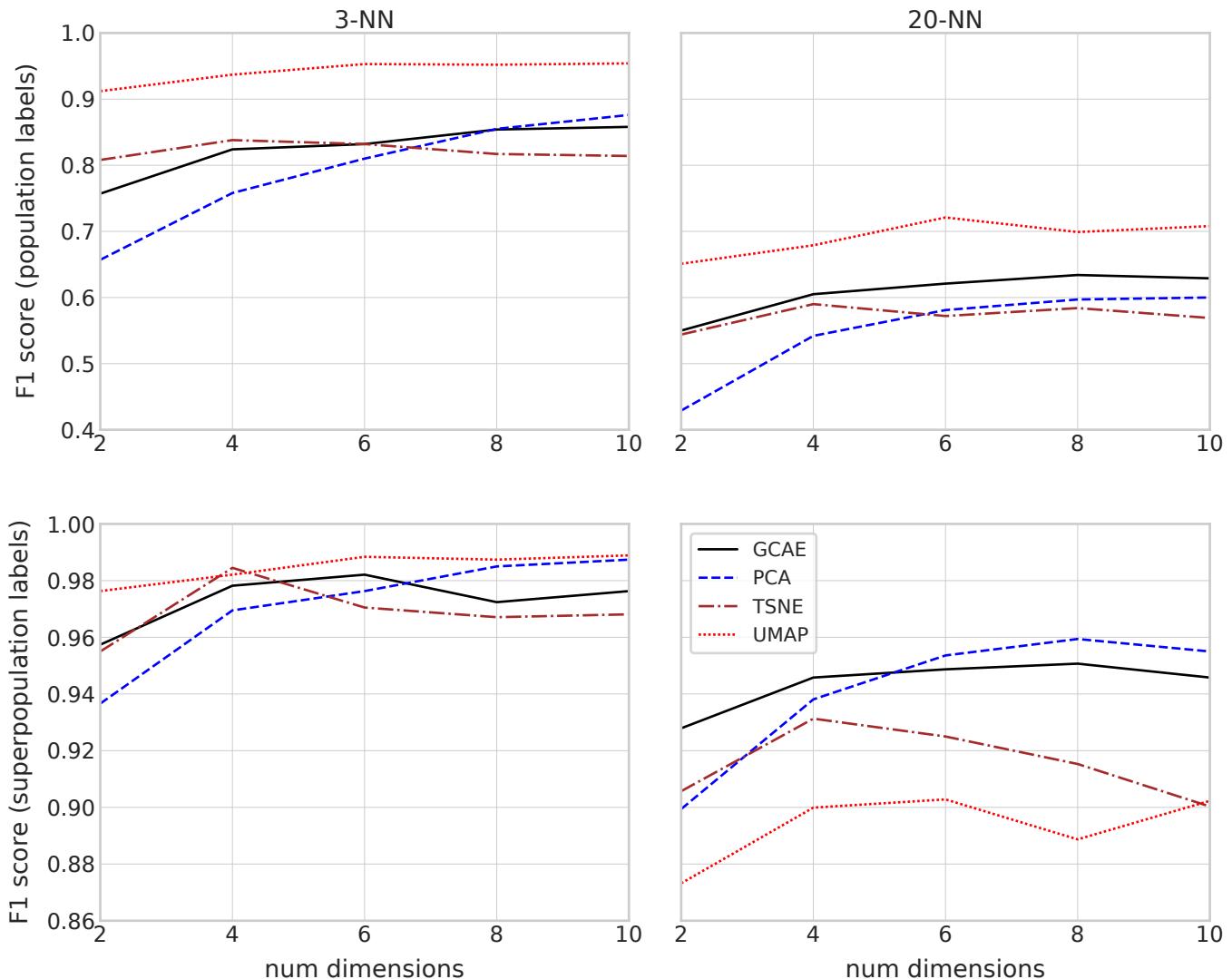


Figure 4 F1 scores for 3-NN and 20-NN classification models based on dimensionality reduction of GCAE, PCA, t-SNE and UMAP using 2-10 dimensions. Top plots show results for using populations as labels in the classification, and bottom ones for using superpopulations as labels.

1 man Origins data set with 5 clusters using ADMIXTURE and
2 GCAE, with the order of the clusters adjusted to be analogous
3 for comparison. Both models result in distinct American, Oceanian
4 and African clusters, with the latter also showing a very similar
5 pattern of the blue component. The ADMIXTURE results further
6 show the East Asian and European superpopulations as largely
7 distinct clusters, whereas GCAE reveals these as composite.

8 For European populations, the most prominent components
9 in the GCAE clustering are blue, which is mainly present in the
10 Middle East for both methods, and red. A possible interpretation
11 is that the red cluster signifies the genetic component of herders
12 that migrated to Europe from the Pontic–Caspian steppe around
13 4.5 kyr ago (Nielsen *et al.* 2017; Haak *et al.* 2015). This would be
14 consistent with a presence in most European populations, with the
15 exception of Sardinians (Lazaridis *et al.* 2014), as well as in South
16 Asia, with particular prominence in e.g. the Kalash (Lamnidis
17 *et al.* 2018; Pathak *et al.* 2018). This component also appears in East
18 Asian populations, which also include American and Oceanian
19 ancestry.

20 Results for the comparison of output data from GCAE to the
21 corresponding true genotypes are shown in Figure 7. The left
22 plot shows that the GCAE-generated genotypes follow the true
23 distribution of allele frequencies in the population. A tendency to
24 underestimate the presence of the derived allele for low-frequency
25 sites, particularly for values below 0.1 in the true data, is however
26 visible, indicating that rare variation is more difficult for GCAE
27 to capture. Another trend visible in the middle of the frequency
28 spectrum is that output genotypes tend to have a higher presence
29 of the derived allele where the true frequency is below 0.5 and a
30 corresponding underestimation of the derived allele where it is in
31 majority in the true data. This is likely an effect of the balancing
32 performed on the loss function to handle uneven classes, and may
33 be mitigated by finer tuning of the β parameter.

34 The right plot shows that output genotypes do show a pattern of
35 decay of LD with distance along the chromosome that reflects that
36 of the true data, with very similar correlation values for the lowest
37 SNP distances. Although the correlation between sites tends to
38 be higher in the GCAE-generated data for longer distances, with
39 a less smooth decay curve, the results suggest that GCAE is able
40 to define the coding into the latent space in a way that takes local
41 spatial structure into account.

42 DISCUSSION

43 One approach to evaluation of dimensionality reduction involves
44 assessing the correspondence between transformed data and its
45 geographical sampling location. Quantitative studies have shown
46 that geographical effects in the form of migration and the impact
47 of physical distance on gene flow play a role in creating population
48 structure (Wang *et al.* 2012). For PCA, striking similarities to geog-
49 raphy have mainly been reported from limited geographical areas
50 such as Europe (Novembre *et al.* 2008; Lao *et al.* 2008), with world-
51 wide cohorts generally resulting in a less resolved V-shape similar
52 to that shown in Figure 3 B (Jakobsson *et al.* 2008; Biswas *et al.* 2009).
53 As previously mentioned, t-SNE and UMAP tend to focus on lo-
54 cal relationships, although Diaz-Papkovich *et al.* (2021) discusses
55 that for UMAP, careful filtering of the data can cause geographical
56 features to be highlighted more, but again, mainly when applied
57 to relatively homogeneous data sets. The Human Origins data
58 set considered here represents worldwide genetic variation, and
59 the visualization results show that GCAE displays robustness to
60 this high degree of diversity, yielding a representation that reflects
61 global geographical patterns.

62 The interpretation of dimensionality reduction results, partic-
63 ularly the inference of the underlying processes behind observed
64 structure, is however not always straightforward. The clustering
65 of samples in reduced-dimensional space can reflect characteristics
66 at various scales in the data, ranging from the presence of a par-
67 ticular variant to continental ancestry. In Novembre and Stephens
68 (2008) and François *et al.* (2010), for example, the effects of past
69 migration and expansion events on PCA is discussed, and how the
70 assumptions of linearity and orthogonality of the model can result
71 in counter-intuitive patterns in PC-space. The effects of attributes
72 of the data such as LD and so-called “informative missingness”,
73 due to e.g. different sequencing panels or the higher uncertainty
74 associated with heterozygote calls, on dimensionality reduction
75 are also extensively discussed in Patterson *et al.* (2006).

76 The F1 score of a classification model based on the dimensionality
77 reduction is not a simple metric for which the method with the
78 highest score is the most correct. Nonetheless, the results indicate
79 relevant, systematic differences between the models. With respect
80 to the performance metric considered, PCA and GCAE both seem
81 to be able to make more efficient use of increasing numbers of
82 dimensions in the latent space than t-SNE and UMAP. Further, the
83 three nonlinear methods tend to reveal more fine-scale patterns
84 and yield a more resolved representation than PCA, at least for
85 lower number of dimensions. A key difference, however, is that
86 GCAE takes a more global approach that preserves the meaning
87 of distances between clusters to a larger extent than the neighbor
88 graph-based methods.

89 Other deep-learning based methods that generate genotypes
90 have shown varying capabilities to preserve spatial properties in
91 terms of decay of LD with distance along the chromosome. The
92 variational autoencoder of Battey *et al.* (2020) failed to reproduce
93 LD decay, while both the generative adversarial network (GAN)
94 and RBM models of Yelmen *et al.* (2021) captured LD patterns
95 well, with correlation at larger distances in particular preserved
96 to a higher extent than for GCAE. An important methodological
97 difference between GCAE and these other models, and also PCA,
98 t-SNE and UMAP, is that they do not take the sequential nature
99 of genotype data into account. Rather, they treat every SNP as
100 an independent variable, whereas convolution treats relationships
101 between nearby sites differently than any arbitrary pair of variants.
102 We show that convolutional architectures can be used to capture
103 local spatial patterns, and believe that additional convolutional
104 and max-pooling layers can improve LD accuracy over longer
105 distances.

106 Our results demonstrate that the use of convolution is feasible
107 for genotype data in spite of its fundamental differences to images,
108 which such networks are typically applied to. Genotype data is
109 position-dependent, with a unique meaning to every dimension.
110 Pixels in images, in contrast, typically represent information that
111 is translation invariant. Our experiments indicated that the incor-
112 poration of positional information in the form of marker-specific
113 variables that the model can optimize during the training process
114 improved performance. We therefore suggest this as a means of
115 allowing the model to represent some of the global information
116 that is lost with convolution. An additional, more explicit, method
117 to include sequential information would be to include e.g. genetic
118 distance or position as part of the input data, as done in Chan *et al.*
119 (2018); Adrián *et al.* (2020), which we leave for future work.

120 We also note that in order to obtain a fair comparison, we have
121 performed filtering of the SNP set in terms of MAF and LD ac-
122 cording to standard protocols for PCA even though these steps are
123 not necessarily required for GCAE, in which such patterns can be

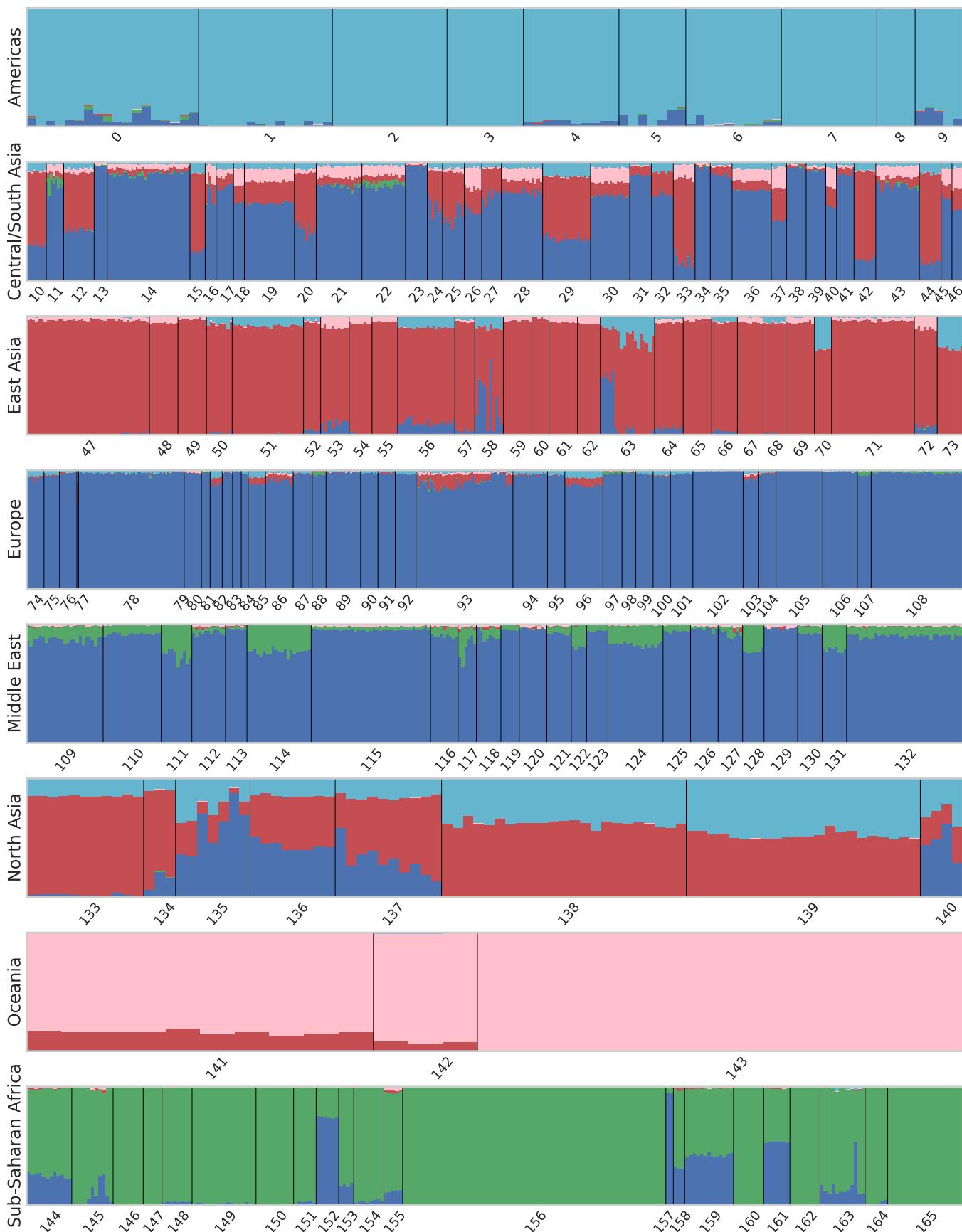


Figure 5 Genetic clustering results with $k = 5$ clusters using ADMIXTURE. Each bar represents a sample from the Human Origins data set, with colors indicating the proportional assignment into k clusters for that sample. Samples are ordered by population and superpopulation, with numbering according to the legend in Figure 1.

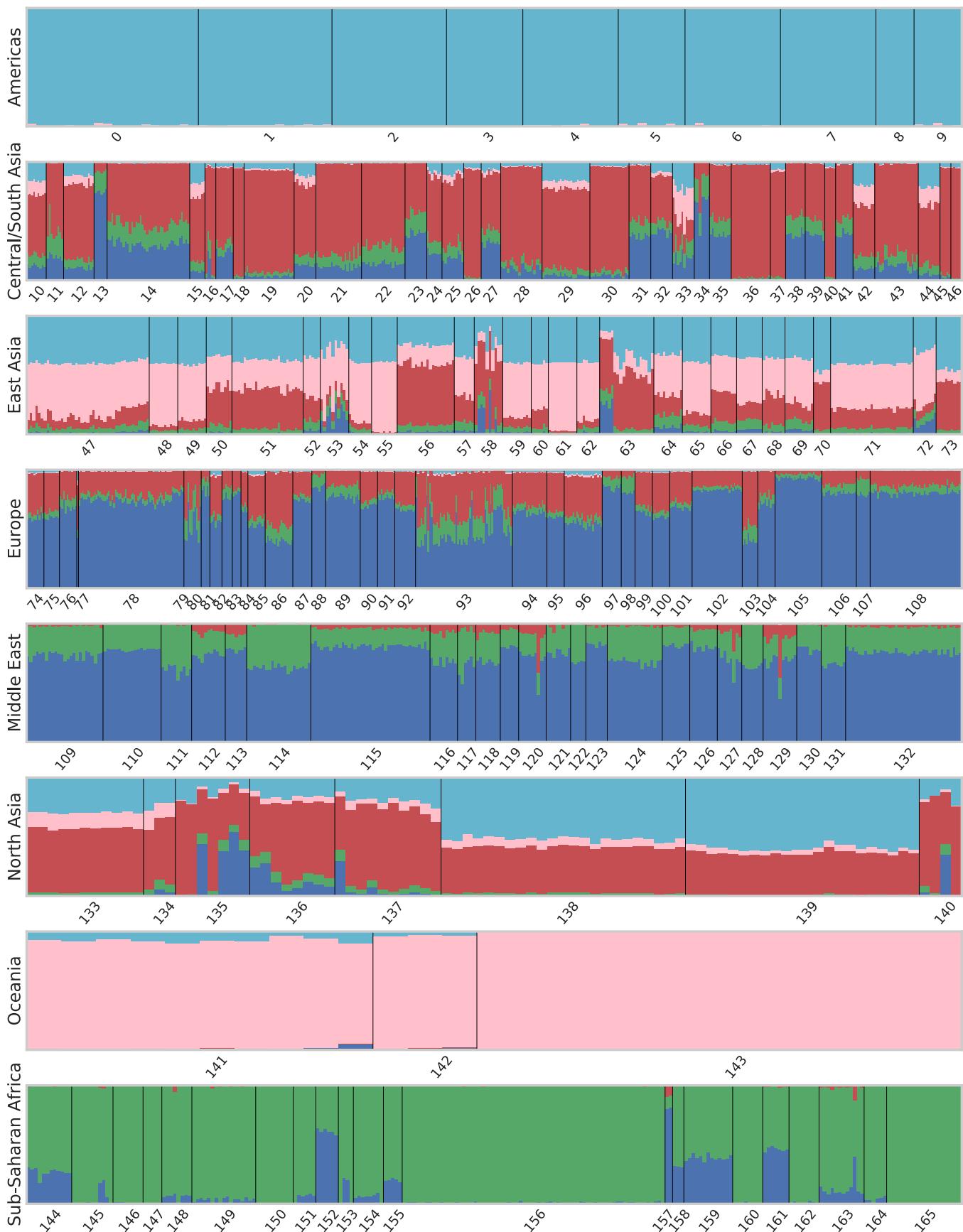


Figure 6 Genetic clustering results with $k = 5$ clusters using GCAE. Each bar represents a sample from the Human Origins data set, with colors indicating the proportional assignment into k clusters for that sample. Samples are ordered by population and superpopulation, with numbering according to the legend in Figure 1.

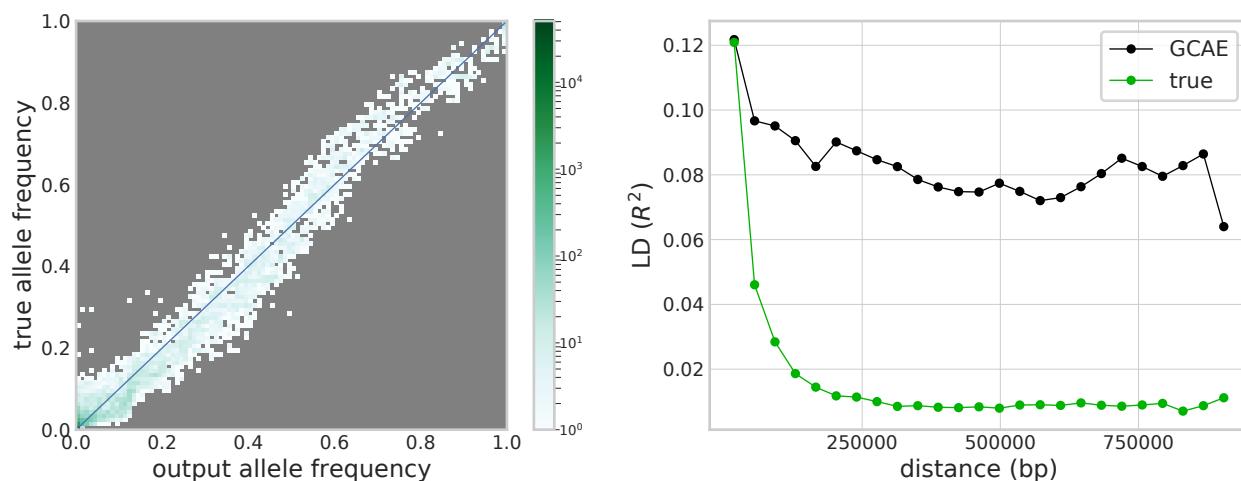


Figure 7 Comparison of output from GCAE to the corresponding true genotypes for the 1000 Genomes chromosome 22 data set. Left: Site frequency spectrum of true genotypes vs. output genotypes, with color indicating number of sites in 100 allele frequency bins. Right: Decay of LD with distance along the chromosome for a subset of samples and SNPs, displaying the mean value of R^2 for pairs of sites in 25 distance bins.

1 learned during the training process. The GCAE architecture also
2 includes a representation of missing genotypes, unlike the other
3 models. In practice, missing data is often handled by either im-
4 putation with the empirical mean and/or filtering to remove sites
5 with high missingness. In this sense, GCAE can present a more
6 robust alternative that is more suitable for low-coverage samples
7 such as ancient DNA, requiring less filtering and allowing more of
8 the data to be retained for analysis.

9 Many commonly used methods for genetic clustering such as
10 TreeMix (Pickrell and Pritchard 2012), ADMIXTUREGRAPH (Lep-
11 päälä *et al.* 2017) and GLOBETROTTER (Hellenthal *et al.* 2014) are,
12 unlike the dimensionality reduction methods discussed, model-
13 based. STRUCTURE (Pritchard *et al.* 2000), for example, represents
14 LD and includes explicit modeling of admixture blocks and the
15 transitions between them. The model assumes the existence of a
16 set of differentiated ancestral populations, and that the sample is
17 a result of relatively recent mixing of these. ADMIXTURE, which
18 we use as a reference method in this work, is based on the same
19 underlying statistical model.

20 GCAE, in contrast, constitutes a more flexible and data-driven
21 approach, which may be an advantage in scenarios where the data
22 does not conform to explicit modeling assumptions. Our results
23 demonstrate that GCAE is able to capture very similar population
24 structure as that found by ADMIXTURE, while also identifying
25 additional characteristics for some populations that are consistent
26 with existing findings in the literature.

27 As previously discussed regarding dimensionality reduction,
28 interpretation and evaluation of genetic clustering is not straight-
29 forward. The correctness of an assignment is not well-defined, and
30 different underlying processes can give rise to similar observed pat-
31 terns (Lawson *et al.* 2018). Evaluation of results requires additional
32 information, such as putting them into the context of methods
33 based on different underlying models, and the analysis of metrics
34 like F and D statistics.

35 A purely data-driven black-box approach such as DL can be dif-
36 ficult to interpret. The features used in the data transformation are
37 unknown and therefore cannot be used for validation of whether
38 certain modeling assumptions hold for the data in question. On
39 the other hand, the alternative methodology of GCAE allows it to
40 capture additional aspects of the data, and therefore provide a use-

ful complement to the toolset used for exploratory data analysis in
population genetics.

Another characteristic of convolutional layers is that they re-
quire less trainable variables than a corresponding fully-connected
layer, leading to reduced computational requirements for training.
Depending on overall network architecture and training strategy,
this may allow for the design of models that are more feasible to
train on large data sets.

When running on CPUs, using 11 cores on UPPMAX, training
of the dimensionality reduction models took between 26.9 and
99.8 hours. Using the GPU on NSC, times ranged between 1.5
and 5.9 hours. The genetic clustering model took 46.6 hours on
CPU, and the model trained on 1000 Genomes data for which
output genotypes were analyzed took 45 minutes on GPU. As a
comparison, PCA took 11 minutes, and t-SNE and UMAP ranged
between 20 minutes and 4 hours on the CPU setup on UPPMAX.

The computational requirements of GCAE are thus greater than
that of the other models, although the use of GPUs can improve
performance significantly. As the purpose of this study is to eval-
uate the applicability of convolutional autoencoders to the chosen
problems, optimization of computational efficiency is considered
out of scope and left for future work.

Our results demonstrate that GCAE can learn features that char-
acterize genotype data in a meaningful way. The minor model
changes required to change the application from dimensionality
reduction to genetic clustering further demonstrate the flexibil-
ity of the method, and future efforts will involve investigating
the application of GCAE to other problems. A simple alterna-
tive application would be imputation of missing genotypes. As the
training procedure is based on reconstructing the input, and since
we already include a representation of missing data, this would
mainly involve finding a suitable number of units to use in the
latent layer. The model can also be used for generation of artificial
genotypes by entering data into the decoder that does not corre-
spond to the encoding of an empirical sample. This can be done
e.g. by perturbing the encoding of actual individuals or selecting
values from a specific part of latent space. If the space is regular
enough, one could use the clustering it has defined to simulate
samples from a particular population or some other characteristic
property learned by the model, e.g. ancient data. We are also

1 currently exploring the use of GCAE in the context of quantitative
2 genetics by incorporating phenotypic information into the model.

3 DATA AVAILABILITY

4 The fully public Affymetrix Human Origins present-day indi-
5 viduals from Lazaridis *et al.* (2016) are available for download
6 from <https://reich.hms.harvard.edu/datasets>. The 1000 Genomes
7 phase 3 data set is available at <https://www.internationalgenome.org/>
8 [data-portal/data-collection/phase-3](https://www.internationalgenome.org/).

9 ACKNOWLEDGMENTS

10 The authors acknowledge the use of computational resources pro-
11 vided by Swedish National Infrastructure for Computing (SNIC)
12 and associated centers under projects Berzelius-2021-30, UPPMAX
13 2020/2-5, SNIC 2019/8-38 and SNIC 2020/5-91. CN also acknowl-
14 edges funding by Formas (grant numbers 2017-00453, 2020-00712).

15 LITERATURE CITED

16 Abadi, M., A. Agarwal, P. Barham, E. Brevdo, Z. Chen, *et al.*, 2015
17 TensorFlow: Large-scale machine learning on heterogeneous
18 systems. Software available from tensorflow.org.
19 Adrión, J. R., J. G. Galloway, and A. D. Kern, 2020 Predicting the
20 Landscape of Recombination Using Deep Learning. Molecular
21 Biology and Evolution 37: 1790–1808.
22 Alanis-Lobato, G., C. V. Cannistraci, A. Eriksson, A. Manica, and
23 T. Ravasi, 2015 Highlighting nonlinear patterns in population
24 genetics datasets. Scientific Reports 5: 8140.
25 Alexander, D., J. Novembre, and K. Lange, 2009 Fast model-based
26 estimation of ancestry in unrelated individuals. Genome Re-
27 search 19: 1655–1664.
28 Alipanahi, B., A. Delong, M. T. Weirauch, and B. J. Frey, 2015
29 Predicting the sequence specificities of dna- and rna-binding
30 proteins by deep learning. Nature Biotechnology 33: 831–838.
31 Auton, A., G. R. Abecasis, D. M. Altshuler, R. M. Durbin, D. R. Bent-
32 ley, *et al.*, 2015 A global reference for human genetic variation.
33 Nature 526: 68–74.
34 Battey, C. J., G. C. Coffing, and A. D. Kern, 2020 Visualizing popu-
35 lation structure with variational autoencoders. bioRxiv .
36 Biswas, S., L. B. Scheinfeldt, and J. M. Akey, 2009 Genome-wide
37 insights into the patterns and determinants of fine-scale popu-
38 lation structure in humans. The American Journal of Human
39 Genetics 84: 641 – 650.
40 Brechtmann, F., C. Mertes, A. Matusevičiūtė, V. A. Yépez, Žiga
41 Avsec, *et al.*, 2018 Outrider: A statistical method for detecting
42 aberrantly expressed genes in rna sequencing data. The Ameri-
43 can Journal of Human Genetics 103: 907 – 917.
44 Chan, J., V. Perrone, J. P. Spence, P. A. Jenkins, S. Mathieson,
45 *et al.*, 2018 A likelihood-free inference framework for popu-
46 lation genetic data using exchangeable neural networks. Ad-
47 vances in neural information processing systems 31: 8594–8605,
48 33244210[pmid].
49 Chang, C. C., C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell,
50 *et al.*, 2015 Second-generation PLINK: rising to the challenge of
51 larger and richer datasets. GigaScience 4, s13742-015-0047-8.
52 Chen, J. and X. Shi, 2019 Sparse convolutional denoising autoem-
53 coders for genotype imputation. Genes 10: 652, 31466333[pmid].
54 Cheng, J., T. Y. D. Nguyen, K. J. Cygan, M. H. Çelik, W. G. Fair-
55 brother, *et al.*, 2019 Mmsplice: modular modeling improves the
56 predictions of genetic variant effects on splicing. Genome Biol-
57 ogy 20: 48.

Cui, Y., M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, 2019 Class-
59 balanced loss based on effective number of samples.
60 Danecek, P., J. K. Bonfield, J. Liddle, J. Marshall, V. Ohan, *et al.*, 2021
61 Twelve years of samtools and bcftools. GigaScience 10: giab008,
62 33590861[pmid].
63 Diaz-Papkovich, A., L. Anderson-Trocme, and S. Gravel, 2021
64 A review of umap in population genetics. Journal of Human
65 Genetics 66: 85–91.
66 Ding, J., A. Condon, and S. P. Shah, 2018 Interpretable di-
67 mensionality reduction of single cell transcriptome data with
68 deep generative models. Nature communications 9: 2002–2002,
69 29784946[pmid].
70 Eraslan, G., Z. Avsec, J. Gagneur, and F. J. Theis, 2019a Deep learn-
71 ing: new computational modelling techniques for genomics.
72 Nature Reviews Genetics 20: 389–403.
73 Eraslan, G., L. M. Simon, M. Mircea, N. S. Mueller, and F. J. Theis,
74 2019b Single-cell rna-seq denoising using a deep count autoen-
75 coder. Nature Communications 10: 390.
76 François, O., M. Currat, N. Ray, E. Han, L. Excoffier, *et al.*, 2010
77 Principal Component Analysis under Population Genetic Mod-
78 els of Range Expansion and Admixture. Molecular Biology and
79 Evolution 27: 1257–1268.
80 Gaspar, H. A. and G. Breen, 2019 Probabilistic ancestry maps:
81 a method to assess and visualize population substructures in
82 genetics. BMC Bioinformatics 20: 116.
83 Goodfellow, I., Y. Bengio, and A. Courville, 2016 Deep Learning.
84 MIT Press, <http://www.deeplearningbook.org>.
85 Haak, W., I. Lazaridis, N. Patterson, N. Rohland, S. Mallick, *et al.*,
86 2015 Massive migration from the steppe was a source for indo-
87 european languages in europe. Nature 522: 207–211.
88 He, K., X. Zhang, S. Ren, and J. Sun, 2015 Deep residual learning
89 for image recognition. CoRR <abs/1512.03385>.
90 Hellenthal, G., G. B. J. Busby, G. Band, J. F. Wilson, C. Capelli, *et al.*,
91 2014 A genetic atlas of human admixture history. Science 343:
92 747–751.
93 Hinton, G. E. and R. R. Salakhutdinov, 2006 Reducing the dimen-
94 sionality of data with neural networks. Science 313: 504–507.
95 Jakobsson, M., S. W. Scholz, P. Scheet, J. R. Gibbs, J. M. VanLiere,
96 *et al.*, 2008 Genotype, haplotype and copy-number variation in
97 worldwide human populations. Nature 451: 998–1003.
98 Kingma, D. P. and J. Ba, 2014 Adam: A method for stochastic
99 optimization.
100 Koh, P. W., E. Pierson, and A. Kundaje, 2017 Denoising genome-
101 wide histone ChIP-seq with convolutional neural networks.
102 Bioinformatics 33: i225–i233.
103 Kramer, M. A., 1991 Nonlinear principal component analysis using
104 autoassociative neural networks. AIChE Journal 37: 233–243.
105 Lamnidis, T. C., K. Majander, C. Jeong, E. Salmela, A. Wessman,
106 *et al.*, 2018 Ancient fennoscandian genomes reveal origin and
107 spread of siberian ancestry in europe. Nature Communications
108 9: 5018.
109 Lao, O., T. T. Lu, M. Nothnagel, O. Junge, S. Freitag-Wolf, *et al.*,
110 2008 Correlation between genetic and geographic structure in
111 europe. Current Biology 18: 1241 – 1248.
112 Lawson, D. J., L. van Dorp, and D. Falush, 2018 A tutorial on how
113 not to over-interpret structure and admixture bar plots. Nature
114 Communications 9: 3258.
115 Lazaridis, I., D. Nadel, G. Rollefson, D. C. Merrett, N. Rohland,
116 *et al.*, 2016 Genomic insights into the origin of farming in the
117 ancient near east. Nature 536: 419–424.
118 Lazaridis, I., N. Patterson, A. Mitnik, G. Renaud, S. Mallick, *et al.*,
119 2014 Ancient human genomes suggest three ancestral popula-

- tions for present-day europeans. *Nature* **513**: 409–413.
- LeCun, Y., Y. Bengio, and G. Hinton, 2015 Deep learning. *Nature* **521**: 436–444.
- Leppälä, K., S. V. Nielsen, and T. Mailund, 2017 admixturegraph: an R package for admixture graph manipulation and fitting. *Bioinformatics* **33**: 1738–1740.
- Libbrecht, M. W. and W. S. Noble, 2015 Machine learning applications in genetics and genomics. *Nature Reviews Genetics* **16**: 321–332.
- Ma, S. and G. Shi, 2020 On rare variants in principal component analysis of population stratification. *BMC Genetics* **21**: 34.
- McInnes, L., J. Healy, and J. Melville, 2020 Umap: Uniform manifold approximation and projection for dimension reduction.
- Miles, A., pyup.io bot, M. R., P. Ralph, N. Harding, *et al.*, 2021 cgh/scikit-allel: v1.3.3.
- Nawy, T., 2018 Variants from the deep. *Nature Methods* **15**: 861–861.
- Nielsen, R., J. M. Akey, M. Jakobsson, J. K. Pritchard, S. Tishkoff, *et al.*, 2017 Tracing the peopling of the world through genomics. *Nature* **541**: 302–310.
- Novembre, J., T. Johnson, K. Bryc, Z. Kutalik, A. R. Boyko, *et al.*, 2008 Genes mirror geography within europe. *Nature* **456**: 98–101.
- Novembre, J. and M. Stephens, 2008 Interpreting principal component analyses of spatial population genetic variation. *Nature Genetics* **40**: 646–649.
- Pathak, A. K., A. Kadian, A. Kushniarevich, F. Montinaro, M. Mondal, *et al.*, 2018 The genetic ancestry of modern indus valley populations from northwest india. *The American Journal of Human Genetics* **103**: 918 – 929.
- Patterson, N., A. L. Price, and D. Reich, 2006 Population structure and eigenanalysis. *PLOS Genetics* **2**: 1–20.
- Pearson, K., 1901 On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**: 559–572.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, *et al.*, 2011 Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**: 2825–2830.
- Pickrell, J. K. and J. K. Pritchard, 2012 Inference of population splits and mixtures from genome-wide allele frequency data. *PLOS Genetics* **8**: 1–17.
- Poplin, R., P.-C. Chang, D. Alexander, S. Schwartz, T. Colthurst, *et al.*, 2018 A universal snp and small-indel variant caller using deep neural networks. *Nature Biotechnology* **36**: 983–987.
- Price, A. L., N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, *et al.*, 2006 Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics* **38**: 904–909.
- Pritchard, J. K., M. Stephens, and P. Donnelly, 2000 Inference of population structure using multilocus genotype data. *Genetics* **155**: 945–959, 10835412[pmid].
- Scholz, M., F. Kaplan, C. L. Guy, J. Kopka, and J. Selbig, 2005 Non-linear PCA: a missing data approach. *Bioinformatics* **21**: 3887–3895.
- Schrider, D. R. and A. D. Kern, 2018 Supervised machine learning for population genetics: A new paradigm. *Trends in genetics : TIG* **34**: 301–312, 29331490[pmid].
- Shaun Purcell, Christopher Chang, 2020 Plink 1.9. www.cog-genomics.org/plink/1.9/.
- Sun, Y. V. and S. L. R. Kardia, 2008 Imputing missing genotypic data of single-nucleotide polymorphisms using neural networks. *European Journal of Human Genetics* **16**: 487–495.
- Talwar, D., A. Mongia, D. Sengupta, and A. Majumdar, 2018 Autompute: Autoencoder based imputation of single-cell rna-seq data. *Scientific Reports* **8**: 16329.
- Tian, C., P. K. Gregersen, and M. F. Seldin, 2008 Accounting for ancestry: population substructure and genome-wide association studies. *Human Molecular Genetics* **17**: R143–R150.
- Ulyanov, D., 2016 Multicore-tsne. <https://github.com/DmitryUlyanov/Multicore-TSNE>.
- van der Maaten, L. and G. Hinton, 2008 Visualizing data using t-sne. *Journal of Machine Learning Research* **9**: 2579–2605.
- Wang, C., S. Zöllner, and N. A. Rosenberg, 2012 A quantitative comparison of the similarity between genes and geography in worldwide human populations. *PLOS Genetics* **8**: 1–16.
- Xu, C. and S. A. Jackson, 2019 Machine learning and complex biological data. *Genome Biology* **20**: 76.
- Yelmen, B., A. Decelle, L. Ongaro, D. Marnetto, C. Tallec, *et al.*, 2021 Creating artificial human genomes using generative neural networks. *PLOS Genetics* **17**: 1–22.
- Zou, F., S. Lee, M. R. Knowles, and F. A. Wright, 2010 Quantification of population structure using correlated snps by shrinkage principal components. *Human heredity* **70**: 9–22, 20413978[pmid].
- Zou, J., M. Huss, A. Abid, P. Mohammadi, A. Torkamani, *et al.*, 2019 A primer on deep learning in genomics. *Nature Genetics* **51**: 12–18.