

# 1 **Harnessing deep learning for population genetic inference**

2 Xin Huang<sup>1,2,†</sup>, Aigerim Rymbekova<sup>1,2</sup>, Olga Dolgova<sup>3</sup>, Oscar Lao<sup>4,†</sup> & Martin Kuhlwilm<sup>1,2,†</sup>

3

4 <sup>1</sup>Department of Evolutionary Anthropology, University of Vienna, Vienna, Austria

5 <sup>2</sup>Human Evolution and Archaeological Sciences (HEAS), University of Vienna, Vienna, Austria

6 <sup>3</sup>Integrative Genomics Lab, CIC bioGUNE - Centro de Investigación Cooperativa en Biociencias,  
7 Derio, Biscaya, Spain

8 <sup>4</sup>Institute of Evolutionary Biology, CSIC–Universitat Pompeu Fabra, Barcelona, Spain

9 †e-mail: xin.huang@univie.ac.at, oscar.lao@ibe.upf-csic.es, [martin.kuhlwilm@univie.ac.at](mailto:martin.kuhlwilm@univie.ac.at)

10

## 11 **Abstract**

12 In population genetics, the emergence of large-scale genomic data for various species and  
13 populations has provided new opportunities to understand the evolutionary forces that drive  
14 genetic diversity using statistical inference. However, the era of population genomics presents  
15 new challenges in analysing the massive amounts of variants and genomes. Deep learning has  
16 demonstrated state-of-the-art performance for numerous applications involving large-scale data  
17 Recently, deep learning approaches have started gaining popularity in population genetics and  
18 been applied in various population genetics problems, including identifying population structure,  
19 inferring demographic history and investigating natural selection, due to the advent of massive  
20 genomic datasets, powerful computational hardware and complex deep learning architectures.  
21 Here, we introduce common deep learning architectures and provide comprehensive guidelines  
22 for implementing deep learning models for population genetic inference. We also discuss current  
23 challenges and future directions for applying deep learning in population genetics, focusing on  
24 efficiency, robustness and interpretability.

25

## 26 **Introduction**

27 Population genetics is a more than century-old discipline that harnesses genetic variation within  
28 and between populations to explore evolutionary processes or forces such as mutation,  
29 recombination, natural selection, genetic drift and gene flow<sup>1-3</sup>. The main goal of population  
30 genetics is to investigate how different evolutionary forces and their interactions shape the  
31 population dynamics of genetic variants with a given demographic history, spatial structure and  
32 mating system<sup>1,2,4</sup>. These evolutionary forces can be mathematically modelled as evolutionary  
33 parameters, and the dynamics of genetic variants can be represented by changes in allele

34 frequencies or times to the most recent common ancestor<sup>5, 6</sup>. Therefore, population genetic  
35 inference refers to estimating evolutionary parameters or identifying genetic patterns produced  
36 by evolutionary forces in populations using different approaches. Typical tasks in population  
37 genetic inference include inferring demographic models, identifying population structure, studying  
38 spatial genetic patterns, estimating mutation or recombination rates, detecting signatures and  
39 quantifying strengths of natural selection, as well as investigating admixture or introgression  
40 events and their footprints in genomes, such as ancestry proportions and introgressed fragments  
41 (Table 1). In addition, simulating genomic data is an important approach utilized in population  
42 genetic inference (Fig. 1).

43 With the advent of high-throughput sequencing technologies, genomic data are now  
44 driving a revolution in population genetics. The paradigm of studying population genetics has  
45 been shifted from using small-scale protein electrophoretic variation or molecular markers to  
46 utilizing large-scale genome-wide single nucleotide polymorphisms or structural variations<sup>7-9</sup>. For  
47 instance, in human population genetics, the 1000 Genomes Project<sup>10</sup> provides high-quality  
48 genotype data for ~90 million variants from 26 modern human populations across five continents  
49 with considerable sample size. New datasets for ancient human genomes and non-human  
50 species are constantly becoming available, such as the Allen Ancient DNA Resource<sup>11</sup> and the  
51 1001 Genomes Project<sup>12</sup>. Other datasets with a focus on functional genomics for identifying  
52 causes of diseases, for example, the UK Biobank<sup>13</sup> and the China Kadoorie Biobank<sup>14</sup>, could be  
53 also considered to be utilized for population genetic inference. However, the capacities of humans  
54 and traditional computational approaches to integrate and interpret the massive amounts of  
55 variants and genomes from different populations and species are often exceeded, raising new  
56 challenges in the genomics age.

57 To tackle these challenges, machine learning has provided various successful  
58 computational strategies for population genetics<sup>15</sup>, including hidden Markov models and principal  
59 component analysis. A promising machine learning approach is deep learning, which uses  
60 artificial neural networks (ANNs), consisting of multiple layers of artificial neurons that perform  
61 mathematical operations and form a hierarchical representation of the data to generate  
62 predictions without the need for mathematical modelling from domain-specific knowledge. Thanks  
63 to the data explosion, hardware revolution and model innovation during the past two decades,  
64 deep learning has achieved state-of-the-art performance in various machine learning tasks, such  
65 as computer vision and natural language processing<sup>16</sup>. The astonishing success of deep learning  
66 in automatizing complex pattern recognition tasks has fuelled the translation of deep learning  
67 approaches in fields where large volumes of data exist but no analytical solutions for the questions

68 are available. In genomics, deep learning has already gained popularity for its potential to provide  
69 innovative solutions for diverse problems such as pathogenic variant diagnosis<sup>17</sup> and genome-  
70 wide association studies<sup>18</sup>.

71 Here, we summarize the recent progress in population genetics with deep learning. We  
72 first briefly review traditional approaches for population genetic inference, followed by machine  
73 learning-based techniques in population genetics and classical ANN architectures, including feed-  
74 forward neural networks (FNNs), convolutional neural networks (CNNs), recurrent neural  
75 networks (RNNs) and graph neural networks (GNNs). We also introduce deep generative models  
76 (DGMs), including variational autoencoders (VAEs) and generative adversarial networks (GANs).  
77 We then present novel and promising deep learning models, including transformers<sup>19</sup> and  
78 diffusion models<sup>20</sup>. Finally, we provide guidelines for implementing deep learning-based tools and  
79 discuss current issues and future directions for deep learning in population genetics.

80

## 81 **Traditional population genetic inference**

82 Traditionally, population genetics uses deterministic and stochastic models<sup>21</sup> (Fig. 1).  
83 Deterministic models assume that the dynamics of genetic variants are not affected by random  
84 fluctuations in evolutionary forces, whereas stochastic models are more realistic, as they explicitly  
85 model random effects in evolutionary processes<sup>6, 21</sup>. The Wright–Fisher model — which assumes  
86 a population of a constant size, random mating and non-overlapping generations in the absence  
87 of the effects of mutation, gene flow and natural selection — is a popular stochastic model that  
88 describes allele frequency changes within populations. This model can be understood as a  
89 sequential process using a discrete-time Markov chain from probability theory<sup>22</sup>. Furthermore, it  
90 can be approximated with a continuous Markov process by the diffusion theory, which models the  
91 dynamics of genetic variants over time in populations by approximating their trajectories with  
92 diffusion processes<sup>23</sup>. Another powerful approach for modelling evolution is the coalescent theory  
93 and its extensions<sup>24, 25</sup>, which model the dynamics of genetic variants in a sample by tracing how  
94 they originated from a common ancestor backwards in time. Once a probabilistic model is  
95 established, statistical inference approaches, including maximum likelihood estimation<sup>26</sup>, Markov  
96 chain Monte Carlo<sup>27</sup> and approximate Bayesian computation<sup>28</sup>, can be evaluated with simulated  
97 data and applied to empirical data for parameter estimation (reviewed in ref. 29).

98 Traditional approaches offer advantages, including interpretability under population  
99 genetics theory and the possibility for performance improvement by increasing model complexity.  
100 However, they also have several disadvantages. First, population genetics models often simplify

101 complex biological processes for computational feasibility<sup>5</sup>. For example, genomic data is usually  
102 assumed neutral when inferring demographic models; however, natural selection may mimic  
103 genetic patterns produced by demographic events<sup>30</sup>. Second, traditional tools may not work well  
104 on high-dimensional and large-scale datasets. For example, scalability becomes a challenge for  
105 methods utilizing extended haplotype homozygosity to efficiently detect signals of natural  
106 selection in datasets containing millions of genomes<sup>31–33</sup>. Third, traditional approaches may lack  
107 versatility, as their performance can vary dependent on factors such as sample sizes, underlying  
108 evolutionary models, and phased or unphased data. For example, tools designed for human  
109 demographic history with large-scale data may not perform well when analysing data with small  
110 sample sizes or from non-human species<sup>34</sup>. To address these limitations, deep learning can  
111 provide innovative solutions that differ from novel methodology based on traditional approaches<sup>35</sup>.  
112 <sup>36</sup>.

113

## 114 **What is deep learning?**

115 The primary goal of machine learning is to develop algorithms that can automatically extract  
116 information from data and find solutions for specific tasks<sup>37</sup>. Supervised learning, unsupervised  
117 learning and reinforcement learning are three learning paradigms of machine learning with some  
118 intermediate categories, such as self-supervised learning and semi-supervised learning<sup>37–39</sup>.  
119 Supervised learning constitutes a machine learning paradigm that trains machine learning models  
120 by using data with known labels provided by humans, whereas in unsupervised learning, machine  
121 learning models are trained using data without labels. Reinforcement learning is a paradigm by  
122 which machine learning models are trained using actions with rewards. In self-supervised  
123 learning, machine learning models are trained by first learning representation from unlabelled  
124 data, which are then used to automatically generate labels for downstream supervised learning  
125 tasks. By contrast, semi-supervised learning constitutes first training on a small, labelled dataset  
126 with supervised learning and then training on a large, unlabelled dataset with unsupervised  
127 learning.

128 Many machine learning algorithms are not domain-specific, allowing researchers to  
129 generalize their questions and solve them by choosing existing machine learning algorithms<sup>40</sup>.  
130 For example, identifying population structure from genetic data can be solved by clustering  
131 algorithms without population genetics theory, such as the *K*-means algorithm<sup>41</sup>. Another machine  
132 learning algorithm, hidden Markov models, has various applications in population genetics,  
133 including inferring historical population size changes<sup>42</sup>, archaic introgressed fragments<sup>43</sup> and

134 strengths of natural selection<sup>44</sup>. Traditional machine learning algorithms usually extract features  
135 with domain-specific knowledge from raw data<sup>16</sup>, such as genetic variants. Summary statistics  
136 from population genetics theory can be considered as features extracted from these genetic  
137 variant datasets<sup>45–47</sup>. However, deep learning allows automatic feature extraction from raw data,  
138 as ANNs contain many layers with different possible operations and form deep hierarchical  
139 architectures<sup>16,38</sup> (Fig. 2a). Here, we outline several common ANN architectures that can address  
140 various problems in population genetic inference.

141

## 142 ***Common architectures***

143 FNNs are the most classic architecture of ANNs<sup>48</sup>. In FNNs, information only flows from the input  
144 layer to the output layer through hidden layers<sup>48</sup> (Fig. 2a). Nodes in ANNs are analogous to  
145 neurons, as they receive outputs from the previous layer and usually transform them with  
146 nonlinear activation functions as inputs for the next layer<sup>48</sup> (Fig. 2a). If every node in a hidden or  
147 output layer receives outputs from all nodes in the previous layer as its inputs, this layer is  
148 considered fully connected or dense (Fig. 2b). Fully connected neural networks are not equivalent  
149 to FNNs, as occasionally claimed<sup>36, 49, 50</sup>. The term ‘fully connected’ describes how a layer  
150 connects to its neighbouring layers in an ANN, whereas ‘feed-forward’ describes how information  
151 flows within an ANN. Thus, an FNN is not necessarily fully connected. For example, CNNs are a  
152 type of FNN that is not fully connected<sup>48</sup> (Fig. 2b). Additionally, a fully connected neural network  
153 is not necessarily feed-forward, such as fully connected RNNs<sup>51</sup>. If FNNs only contain fully  
154 connected layers, they are referred to as multilayer perceptrons<sup>52</sup>. Fully connected FNNs or  
155 multilayer perceptrons have illustrated the capability of ANNs in various population genetics  
156 problems, including inferring demographic history<sup>53–56</sup> and population structure<sup>57</sup>, detecting  
157 admixture events<sup>50</sup>, dissecting genomic regions under natural selection<sup>53, 58, 59</sup>, estimating  
158 mutation rates<sup>60, 61</sup> or predicting sample locations from genomes<sup>62</sup>.

159 CNNs are currently the dominant architecture in population genetics (Table 1). CNNs  
160 typically process grid-like inputs, including images, matrices, or tensors. Several data types may  
161 be interpreted in an image-like manner, for example, genotype matrices<sup>63, 64</sup>. A standard CNN  
162 contains several series of convolutional and pooling layers to automatically extract features from  
163 raw data before passing them into fully connected layers (Fig. 2a). In a convolutional layer, a filter  
164 composed by a set of kernels captures information within a small region of the input through a  
165 mathematical operation called convolution. The output from the convolutional layer is usually  
166 transformed by rectified linear unit activation functions or related ones and then often fed into a

167 pooling layer for combining information from adjacent regions (Fig. 2b). Pooling layers make  
168 CNNs invariant to small translations in inputs (Fig. 2b); however, they may be inessential if precise  
169 spatial information is required<sup>48</sup>, as in AlphaGo<sup>65</sup>. CNNs also have been employed to numerous  
170 population genetics tasks, including demographic inference<sup>55, 63</sup>, local ancestry inference<sup>66, 67</sup>,  
171 detection of natural selection<sup>36, 58, 63, 68–77</sup> or introgression<sup>63, 78–81</sup>, and estimation of mutation or  
172 recombination rates as well as dispersal distance<sup>61, 63, 82</sup>. Moreover, fully connected FNN and CNN  
173 architectures can be utilized as components in DGMs<sup>64, 83, 85, 86, 87</sup>. However, treating genotype  
174 matrices as images may cause some problems, as image pixels are row-ordered, whereas  
175 genomes in genotype matrices have no inherent order in population genetics, and evolutionary  
176 parameters are invariant to switching genome order<sup>63, 80</sup>. This problem can be addressed by  
177 making the deep learning architecture insensitive to the order<sup>75, 88</sup> or creating order from sequence  
178 similarity<sup>63, 80</sup>. Hence, architectures should be assessed with domain-specific knowledge when  
179 transferring architectures from machine learning into population genetics.

180 RNNs differ from FNNs because they allow information to flow back to previous layers or  
181 within the same layer<sup>48, 89</sup> (Fig. 2a), thus enabling the tracking of information from previous inputs  
182 using memory. RNNs usually handle sequential inputs, including speech, text, music and  
183 biological sequences. For example, an RNN can accept genetic variants position by position and  
184 determine the output at one position using information from previous positions (Fig. 2c). In  
185 principle, RNNs could memorize all past information; however, in practice, they can only  
186 remember information within a short range depending on the nature of the data and  
187 architectures<sup>48, 89</sup>. Several techniques, such as long short-term memory and gated recurrent units,  
188 have been established to increase memory range and improve RNN performance<sup>48, 89</sup>. RNNs  
189 have been applied to population genetics tasks, including inferring recombination maps and  
190 detecting selective sweeps<sup>90, 91</sup>, outperforming or performing as well as leading methods based  
191 on traditional approaches. However, RNN performance in machine learning has been surpassed  
192 by a novel architecture: transformers<sup>19</sup> (discussed below). Therefore, future studies may consider  
193 transformers rather than RNNs for processing sequential data in population genetics.

194 GNNs utilize graphs as inputs. Graphs are models to represent relationships among  
195 entities using nodes and edges<sup>92</sup>. For instance, ANNs themselves can be visualized with graphs,  
196 where nodes represent operations with results and edges represent parameters (Fig. 2). Two  
197 types of GNNs are widely used. One is graph convolutional networks, which generalize  
198 convolutional layers on grid-like data to graph-structured data<sup>93</sup> (Fig. 2d). The other is graph  
199 attention networks, which replace graph convolutional layers with self-attention layers<sup>94</sup>. Graphs  
200 are common in population genetics. For example, ancestral recombination graphs are used to

201 reconstruct the evolutionary history among genomes with recombination<sup>95</sup>. Haplotype networks  
202 are frequently applied for analysing and visualizing relationships between haplotypes<sup>96</sup>. GNNs  
203 have been used for demographic inference from ancestral recombination graphs<sup>97</sup>, demonstrating  
204 the potential of GNNs in population genetics.

205

## 206 ***Deep generative models***

207 Both supervised and unsupervised learning can apply the architectures above. Unsupervised  
208 learning is challenging because ground truth is absent for determining model performance.  
209 Generative models, an unsupervised learning approach, can learn the intrinsic properties of  
210 training data and create similar data that are not in the training dataset. While non-ANN generative  
211 models such as hidden Markov models have been popular in population genetics for decades,  
212 DGMs are more powerful for many questions<sup>98</sup>. Besides VAEs and GANs, other DGMs include  
213 energy-based models, which map the probability of data into an energy function and link good  
214 predictions with low energy and poor predictions with high energy; normalizing flows, which  
215 gradually approximate the complicated probability distribution of data starting from a simple initial  
216 probability distribution, achieved through a series of invertible and differentiable mathematical  
217 transformations; and autoregressive models, which process sequential data and make predictions  
218 in one step using the predictions from previous steps as inputs<sup>99</sup>. A classic energy-based model  
219 is restricted Boltzmann machines<sup>48</sup> (Fig. 3a), which contain only one visible (input) layer and only  
220 one hidden layer, with connections only between these two layers but not within each layer; this  
221 ANN for unsupervised learning has recently been applied for simulating human genomes<sup>87, 100</sup>.  
222 However, restricted Boltzmann machines have been overshadowed by other DGMs in machine  
223 learning owing to potential limitations, including their lack of training efficiency and scalability to  
224 large datasets, inflexibility in implementing various tasks, and inability to generate high-quality  
225 outputs<sup>101, 102</sup>.

226 VAEs are based on autoencoders, which contain encoders and decoders<sup>48</sup> (Fig. 3b). The  
227 encoder transforms an input into a latent representation termed code, whereas the decoder  
228 attempts to reproduce the original input using the code generated by the encoder<sup>48</sup>. Classic  
229 autoencoders serve as non-linear dimensionality reduction tools; however, they are not  
230 generative models, as they can only reconstruct their inputs and cannot create useful new data  
231 from their codes<sup>48</sup>. VAEs are DGMs that extend autoencoders, typically by minimizing the  
232 distance between the distribution of their codes and a normal distribution<sup>48</sup>, and can synthesize  
233 meaningful new data by randomly sampling a code from the learnt distribution of their codes (Fig.

234 3b). Autoencoders<sup>103</sup> and VAEs<sup>86</sup> are faster and more scalable for estimating ancestry proportions  
235 than traditional approaches such as ADMIXTURE<sup>104</sup> and ChromoPainter<sup>105</sup>; however, they are  
236 slower for identifying population structure compared to dimensionality reduction techniques such  
237 as principle component analysis<sup>83, 84</sup>.

238 GANs are capable of synthesizing higher-quality outputs than VAEs<sup>48, 101, 102</sup>. A typical  
239 GAN contains a generator and a discriminator engaged in a competitive process<sup>48</sup>. The  
240 discriminator tries to distinguish real data from synthetic data generated by the generator, and the  
241 generator aims to create convincing synthetic data using noise from normal distributions, making  
242 it difficult for the discriminator to identify the fakes<sup>48</sup> (Fig. 3c). GANs or GAN-like architectures  
243 (Table 1) have been applied to inferring demographic models<sup>64</sup>, detecting loci under selection<sup>106</sup>,  
244 estimating recombination rate and demographic parameters simultaneously<sup>107</sup>, and creating  
245 individual genomes from real or simulated data<sup>85, 87, 100, 108</sup>. Two issues may limit the practical  
246 applications of GANs. First, training GANs can be unstable due to the alternating parameter  
247 updates between the generator and discriminator during training<sup>48</sup> (Fig. 3c). If the discriminator  
248 fails to learn how to identify synthetic data, the generator will also struggle to learn how to create  
249 realistic synthetic data. Second, GANs may not generate diverse outputs<sup>109</sup>. They might forget  
250 how to produce certain types of data, resulting in mode dropping, or they may only generate a  
251 limited variety of distinct data, leading to mode collapse while still achieving good performance  
252 (Fig. 3d). Outputs from VAEs and GANs can reflect population structure and produce allele  
253 frequency spectra similar to their training data; however, they may not adequately capture  
254 patterns of linkage disequilibrium from the training data<sup>83, 85, 87, 100</sup>.

255

### 256 ***Novel architectures and models***

257 Since deep learning is a rapidly evolving field, novel architectures and DGMs are continually  
258 emerging. One promising architecture is transformers<sup>38</sup>, which are also deep FNNs<sup>110</sup> (Fig. 4a).  
259 They are related to RNNs<sup>111</sup>, and their core components are self-attention layers (Fig. 4b) that  
260 can mimic convolutional layers<sup>112</sup>. Architectures constructed with self-attention layers have  
261 surpassed CNNs and RNNs in different machine learning tasks<sup>113, 114</sup>. Even in biological studies,  
262 transformer-based models can also outperform RNN-based models when predicting mutations  
263 that cause immune escape and affect the fitness of the SARS-CoV-2 virus<sup>115</sup>. Besides self-  
264 attention, many variants of attention mechanisms have been developed and applied in various  
265 transformer-based architectures<sup>110</sup>. The first transformer architecture was developed for  
266 sequence-to-sequence learning with the encoder–decoder architecture<sup>19, 110</sup>. Moreover, the



267 encoder can function independently, similar to BERT (Bidirectional Encoder Representations from  
268 Transformers)<sup>116</sup>, and the decoder can likewise be used in isolation, such as generative pre-  
269 trained transformers (GPT)<sup>117</sup>. Transformers are data-hungry and require substantial  
270 computational memory for training<sup>118, 119</sup>. Therefore, they may not surpass other deep learning  
271 architectures for the same task without sufficient training data<sup>118</sup>. As more large-scale datasets  
272 become available, transformers may show potential for fully utilizing big data in population  
273 genetics. Diffusion models are another promising DGMs that have recently outperformed GANs  
274 in image synthesis<sup>Error! Reference source not found.</sup>. Diffusion models contain a forward diffusion process  
275 and a reverse denoising process<sup>121</sup> (Fig. 4c). During the forward process, inputs are incrementally  
276 added with noise, while an ANN attempts to recover the inputs by removing the noise during the  
277 reverse process<sup>121</sup>.

278 The above architectures can be used alone or in combination. For example, integrating a  
279 VAE and GAN can result in better performance than using a VAE alone for simulating human  
280 genomes<sup>108</sup>. Furthermore, ANNs can be combined with other machine learning or statistical  
281 inference approaches, for example, approximate Bayesian computation with deep learning for  
282 demographic inference from allele frequency spectra<sup>54</sup> or mixture models with ANNs for inferring  
283 the distribution of fitness effects, that is, the proportions of deleterious, neutral and beneficial  
284 mutations across the genome<sup>122</sup>. Hence, there are numerous opportunities for enhancing the  
285 performance of deep learning-based tools and exploring various problems in population genetics  
286 using different architectures and their combinations.

287

## 288 **How to implement deep learning tools?**

### 289 ***Model optimization***

290 The success of applying deep learning-based approaches relies on not only the knowledge of  
291 architecture but also the implementation of deep learning models<sup>48</sup> (Fig. 5). Before implementing  
292 a deep learning model, training data should be collected and may be pre-processed (Box 1).  
293 However, data should be cautiously pre-processed because artifacts might be introduced, making  
294 the deep learning model overfit the training data<sup>75</sup>. In addition, several parameters termed  
295 hyperparameters need to be specified before optimization, as they cannot be optimized  
296 automatically. These hyperparameters are usually related to the architecture itself, such as the  
297 number of neurons in each layer, the number of hidden layers and the form of activation functions.  
298 Hyperparameters are also involved in the training procedure, such as batch size and the number  
299 of epochs (Fig. 5). Different settings of hyperparameters could affect the performance of a deep

300 learning model, and they may be tuned with the validation set by approaches such as grid search  
301 or random search<sup>48</sup> (Fig. 5).

302 A deep learning model is usually optimized by finding the best-fit parameters that minimize  
303 a loss function depending on tasks and data. For example, the mean squared error, which  
304 calculates the average distance between all true values in the training data and their values  
305 predicted by the model, or the root mean squared error are suitable for regression, whereas the  
306 cross-entropy may be desirable for classification. Usually, the loss function can be interpreted as  
307 the negative logarithm of the likelihood function in a deep learning model<sup>123</sup>; therefore, optimizing  
308 a deep learning model by minimizing the loss function can be regarded as a maximum likelihood  
309 approach<sup>123</sup>. One common optimization approach is stochastic gradient descent with  
310 backpropagation<sup>48, 123</sup>. Stochastic gradient descent has a key parameter termed learning rate,  
311 which is also a crucial hyperparameter that controls how much deep learning model parameters  
312 should be updated in each iteration by affecting gradient magnitudes<sup>48, 124</sup>. A large learning rate  
313 may make a deep learning model unable to find the optimal parameters, whereas a small learning  
314 rate may slow down the training procedure<sup>48</sup>. Hence, choosing a suitable learning rate is  
315 challenging. One solution is learning rate schedulers that can change learning rates with a plan  
316 before training. Another solution is adaptive learning rate methods that can vary learning rates  
317 based on gradients during training, such as RMSProp and Adam<sup>125, 126</sup>. The calculation of  
318 gradients becomes infeasible when the loss function lacks differentiability. In such cases,  
319 alternative optimization techniques can be explored, including Monte Carlo approaches<sup>107</sup> and  
320 natural computing algorithms such as simulated annealing approaches<sup>64</sup>. Moreover, deep  
321 learning models can be optimized by treating their parameters as random variables and using the  
322 Bayesian paradigm to estimate their posterior distributions<sup>123</sup>. However, the computational  
323 complexity of Bayesian neural networks is typically higher than that of classic neural networks  
324 (see the tutorials in ref. 127).

325 Several approaches can improve deep learning model performance. Underfitting on  
326 training data should be examined first. Increasing deep learning model complexity can reduce  
327 underfitting, for example, by increasing the number of hidden layers or neurons, or by changing  
328 the form of the activation function. If a complex deep learning model still performs poorly on  
329 training data, adjusting the optimization algorithm can be a solution, for example, using a different  
330 learning rate or optimization algorithm. When the validation loss is small or acceptable, a deep  
331 learning model is optimal, and its performance can be measured on test data (Fig. 5). If a deep  
332 learning does not perform well on test data, it may overfit training data. In this case, the deep  
333 learning model could be improved by reducing the complexity of the deep learning model even if

334 increasing the error in the training dataset. Another strategy may be collecting more training data  
335 or modifying existing training data to increase the size and diversity of training data through data  
336 augmentation<sup>48</sup>. Regularization techniques, including weight decay<sup>123</sup>, dropout<sup>128</sup> and early  
337 stopping<sup>129</sup>, could also reduce overfitting by adding constraints into the machine learning models.  
338 Other techniques, including cross-validation<sup>130</sup> and batch normalization<sup>131</sup>, have regularizing  
339 effects as well<sup>123, 132</sup>, although their primary purposes are not regularization. For a deep learning  
340 model, the best regularization techniques need to be determined by experiments. For instance,  
341 dropout and batch normalization seemed not to work in a VAE implementation for visualizing  
342 population structure<sup>83</sup>. Simple models such as logistic regression or decision trees can be used  
343 as baseline models before applying deep learning models<sup>48</sup>. If the baseline models achieve  
344 acceptable performance, then there may be no need to experiment with deep learning methods,  
345 which are usually computationally intensive and less explainable. If a deep learning model cannot  
346 outperform such baseline models, the implementation should be improved or abandoned.

347

### 348 ***Performance measures***

349 To assess performance of the deep learning model, various measures can be chosen based on  
350 data and tasks<sup>48</sup>. Supervised learning tasks can be categorized into regression, classification, and  
351 structured prediction. For regression, besides the mean squared error or the root mean squared  
352 error, correlation coefficients can evaluate performance, such as the Pearson or Spearman  
353 correlation coefficients. For binary classification, several common measures, including accuracy,  
354 precision, recall and F1 score, are based on the numbers of correct and incorrect predictions.  
355 These can be visualized with confusion matrices, receiver operating characteristic curves or  
356 precision–recall curves. Since different measures quantify different aspects of deep learning  
357 model performance, appropriate measures should be selected according to specific questions.  
358 For example, archaic introgressed fragments are rare in modern human genomes<sup>133</sup>. Precision  
359 and recall may be good choices in this case, because researchers often want to know how many  
360 inferred introgressed fragments are true and how many true introgressed fragments can be  
361 identified overall<sup>34</sup>. These measures can be further extended into multiclass classification and  
362 structured prediction.

363 DGMs are usually assessed by the fidelity and diversity of their outputs<sup>134</sup>; however,  
364 evaluating their performance poses a challenge owing to the computationally intractable nature  
365 of their likelihood functions<sup>135, 136</sup>. Although various measures have been proposed for evaluating  
366 DGMs, a consensus on the best one has yet to be reached<sup>137</sup>. Common measures in image

367 synthesis are unsuitable for population genetic inference because they are specific to ANNs  
368 trained on image datasets<sup>64</sup>. Previous studies have employed several qualitative and quantitative  
369 approaches to evaluate DGMs in population genetics. Qualitative measures can be obtained by  
370 comparing dimensionality reduction results between real data and synthetic data generated by  
371 the deep learning model<sup>83, 87, 100, 138</sup>, assessing summary statistics from population genetics theory  
372 for both real and synthetic data<sup>64, 83, 85, 87, 100</sup>, or contrasting the results from DGMs with those from  
373 tools using traditional approaches<sup>64, 86</sup>. Quantitative measures can be derived by the classification  
374 accuracy of discriminators or other supervised classifiers when distinguishing between real and  
375 synthetic data<sup>64, 138</sup>. Some studies have also utilized the nearest-neighbour adversarial accuracy  
376 to measure overfitting and underfitting<sup>85, 87, 100</sup>, which calculates the average distance between all  
377 data points and their nearest neighbours in both the real and synthetic datasets (details in ref.  
378 100). The resulting score ranges from 0 to 1; a score >0.5 suggests overfitting, whereas a score  
379 <0.5 indicates underfitting. Precision and recall for probability distributions are additionally  
380 applicable to assessing generative models by defining precision as the probability that a random  
381 data point from the probability distribution of the synthetic dataset falls within the support of the  
382 probability distribution of the real dataset, and by defining recall as the probability that a random  
383 data point from the probability distribution of the real dataset falls within the support of the  
384 probability distribution of the synthetic dataset<sup>136, 139</sup>. With these definitions, the performance of  
385 DGMs on different kinds of data, such as images and text, can be evaluated by using the same  
386 metrics<sup>136</sup>. These metrics can be further improved for better assessing and understanding  
387 DGMs<sup>139</sup>. Besides, other recent advancements proposed to evaluate GANs (reviewed in refs. 109  
388 and 134) might also prove useful in population genetics.

389

## 390 **Current issues and future directions**

### 391 ***The chances for deep learning in population genetics***

392 Experimental ANNs for population genetics can be traced back to approximately 30 years ago<sup>140</sup>.  
393 A later study suggested that the performance of ANNs was marginally better than traditional  
394 likelihood-based approaches for inferring the origin of individuals at that moment<sup>141</sup>. The current  
395 success of deep learning is mainly due to the avalanche of big data, the advances of complex  
396 deep learning architectures and the advent of powerful graphics processing units<sup>38, 142</sup>. Applying  
397 deep learning in population genetics has several advantages now. First, deep learning provides  
398 algorithmic approaches for approximating complicated mathematical functions<sup>143, 144</sup>. Hence,

399 deep learning does not rely on likelihood functions derived from population genetics theory and  
400 allows for end-to-end learning<sup>18, 48</sup>, which involves training a machine learning model to predict  
401 final outputs directly from raw inputs, bypassing the need for manually defined features or dividing  
402 the task into multiple steps. By identifying complex non-linear patterns present in the data without  
403 the need for simplified model assumptions, deep learning models have the potential to outperform  
404 traditional approaches. Second, deep learning can effectively handle and extract information from  
405 high-dimensional data<sup>123</sup>, such as large vectors or tensors of genotypes, summary statistics, or  
406 many complex and highly correlated features, whereas traditional approaches usually use one or  
407 a few summary statistics. Third, deep learning can efficiently analyse large-scale data with many  
408 genomes for tasks, such as estimating ancestry proportions in datasets containing hundreds of  
409 thousands of genomes<sup>86, 103</sup>. Fourth, deep learning offers a versatile approach capable of  
410 processing either phased or unphased data with both large and small sample sizes in population  
411 genetics<sup>82</sup>. Researchers then could apply a single tool for analysing different types of data,  
412 because unphased data with small sample sizes are common in non-human species<sup>82</sup>. Fifth, deep  
413 learning can uncover patterns undetectable by traditional approaches. For example, deep  
414 learning models inferred an archaic introgression event from a third archaic human population  
415 before the split of East Asians, South Asians and Oceanians<sup>54</sup>, whereas only Neanderthal and  
416 Denisovan ancestries were previously detected in these populations<sup>145</sup>. Sixth, deep learning  
417 provides efficient methods to synthesize realistic genomes in research without violating privacy<sup>85,</sup>  
418 <sup>87, 100, 108, 138</sup>. Finally, many deep learning approaches have not yet been explored in population  
419 genetic inference. For example, it will be intriguing to compare the performance of diffusion  
420 models with other DGMs for simulating genomes. Similar to text-to-image synthesis<sup>146</sup>, future  
421 studies could explore how to generate artificial genomes from plain text description. Furthermore,  
422 implementing reinforcement learning would be interesting, as it may mimic evolution<sup>147</sup>, design  
423 novel algorithms<sup>148, 149</sup>, or optimize deep learning models with non-differentiable loss functions<sup>150</sup>.  
424 Moreover, it could enable ANNs to operate without training data derived from human knowledge  
425 while surpassing human experts, similar to AlphaGo Zero<sup>151</sup>. Nevertheless, the nature of  
426 population genetic data and problems bring unique obstacles.

### 427 ***The challenges for deep learning in population genetics***

428 The current main issues involve the efficiency, robustness and interpretability of deep learning  
429 models. Training deep learning models can be time-consuming and resource-intensive because  
430 many parameters and hyperparameters need fine-tuning<sup>39, 152</sup>. Trends in machine learning involve  
431 building deep learning models by increasing numbers of layers and parameters<sup>153</sup>, like ResNet,

432 which can comprise 1,000 layers<sup>154</sup> or GPT and other large models, which can contain billions of  
433 parameters<sup>155–158</sup>. Although these complex deep learning models have achieved impressive  
434 performance in machine learning, academic researchers may struggle to train similar models with  
435 limited computational resources. To improve training efficiency, several strategies could be  
436 considered. First, recycling pre-trained models as a starting point may prove beneficial<sup>39</sup>. Several  
437 recently developed deep learning-based tools have provided pre-trained models, enabling the  
438 direct analysis of user-defined data (Table 1). These pre-trained models can reduce the amount  
439 of training data and time while enhancing performance on similar problems through transfer  
440 learning. For example, training a deep learning model to estimate the mutation rate of *de novo*  
441 mutations can begin with another pre-trained deep learning model for rare mutations and  
442 outperform those trained from scratch<sup>61</sup>. Moreover, pre-trained large language models on text can  
443 improve classification tasks on DNA sequences, indicating the potential of reusing pre-trained  
444 models from other fields<sup>159</sup>. Second, developing efficient architectures and algorithms remains an  
445 open-ended challenge (reviewed in ref. 153). For example, network compression can remove  
446 unnecessary parameters and reduce deep learning model sizes<sup>160</sup>. Meta-learning can optimize  
447 hyperparameters, search neural network architectures and improve training efficiency on various  
448 tasks with limited data<sup>161–163</sup>. Third, meta-heuristic approaches may serve as alternative  
449 optimization methods to stochastic gradient descent<sup>164</sup>. Of particular interest is neuroevolution, in  
450 which parameters and hyperparameters of deep learning models are fitted by mimicking  
451 evolution<sup>165, 166</sup>.

452 Building robust deep learning models presents several challenges. First, issues with  
453 training data can have a significant impact on deep learning model performance through various  
454 forms. One form is data shortage, commonly existing in genomic datasets from non-human  
455 species, even for our closest primate relatives<sup>167</sup>. Besides collecting more data from non-human  
456 species, possible solutions include few-shot learning<sup>168</sup>, which enables deep learning models to  
457 be trained with only a few training samples, and zero-shot learning<sup>169</sup>, which allows deep learning  
458 models to make inferences with unseen data. For example, SALAI-Net is a species-agnostic ANN  
459 that can first be trained with sufficient human data and then applied for local ancestry inference  
460 with non-human data<sup>67</sup>. Another challenge is data imbalance, when training data contain uneven  
461 proportions of data from different categories in classification, potentially biasing the model  
462 performance towards the majority category. Data imbalance is frequently observed in genomic  
463 features such as genetic variants under recombination or recent selection<sup>170</sup>. Potential solutions  
464 may be balancing the proportions of different groups in the training data before the start of training  
465 or developing deep learning models that can be trained using imbalanced datasets, like in the

466 case of detecting ghost introgressed fragments with introUNET<sup>80</sup>. Other forms include data  
467 mismatch, whereby the validation or test data have different statistical properties from the training  
468 data, potentially leading to poor model performance on unseen data. Data mismatch can be due  
469 to model misspecification. For example, tools developed under the standard Wright–Fisher model  
470 may not perform well with data from other evolutionary models<sup>97</sup>.

471 Second, supervised learning (Table 1) requires knowing the true evolutionary parameters  
472 in training data, which is usually impossible in population genetics. One solution may be simulated  
473 data from known evolutionary models, as performing realistic simulations is increasingly possible  
474 with more available curated datasets and demographic models<sup>171</sup>. Popular population genetics  
475 simulators, such as ms and msprime<sup>172, 173</sup>, are coalescent approximations to the neutral Wright–  
476 Fisher model, although others like SLiM can evolve forwards in time with natural selection and  
477 implement non-Wright–Fisher models<sup>174</sup>. Hence, tools trained on simulated data may have  
478 different performances on real data<sup>39</sup>. For example, although background selection is common in  
479 real data, it is usually not considered in simulations, because simulating such data is more time-  
480 consuming<sup>175</sup>. However, ignoring background selection can bias demographic inference<sup>176</sup>.  
481 Domain adaptation can help supervised deep learning models achieve good performance even  
482 on training data without background selection or from mis-specified demographic models<sup>177</sup>.  
483 Another solution may be unsupervised learning, utilizing training data without known true  
484 evolutionary models<sup>64, 85, 86</sup>. Additionally, pre-training and self-supervised learning can improve  
485 model robustness<sup>178, 179</sup>.

486 Third, deep learning models and tools should be reproducible. Currently, TensorFlow<sup>180</sup>  
487 and PyTorch<sup>181</sup> dominate deep learning frameworks (Table 1). However, a good deep learning  
488 model in one framework may not perform well in another. Even within the same framework, their  
489 performance may vary due to factors like hardware, random seed, or the framework version<sup>182</sup>.  
490 Therefore, deep learning model development should be documented and reported<sup>183</sup>. Future  
491 studies should consider applying and developing tools for improving reproducibility and  
492 reusability. Open Neural Network Exchange offers a potential solution for interoperability between  
493 different deep learning frameworks. dnadna is a recently developed framework with PyTorch as  
494 backend, aiming to improve the development of deep learning-based tools for reproducibility and  
495 user-friendliness in population genetics<sup>184</sup>.

496 Fourth, deep learning models may be susceptible to adversarial attacks. A recent  
497 technique was developed for creating adversarial mutations on genomes, and even a small  
498 proportion of adversarial mutations could dramatically reduce the deep learning model  
499 performance on local ancestry inference<sup>185</sup>. It remains unknown how such adversarial attacks

500 would affect deep learning in other population genetics tasks. For example, bad actors may slightly  
501 modify their data to fool deep learning tools to construct erroneous evolutionary histories based  
502 on their own interests. Consequently, future studies should investigate adversarial attacks and  
503 defence techniques when applying deep learning<sup>186</sup>.

504 Ultimately, developing interpretable deep learning models should be the goal, as deep  
505 learning models should provide not only accurate predictions but also insightful explanations<sup>39</sup>.  
506 <sup>187</sup>. Reducing the deep learning model complexity may be one strategy to help researchers  
507 interpret deep learning models<sup>75, 103</sup>. Other strategies may include techniques from explainable  
508 artificial intelligence (xAI), which aim to understand the decision procedures of deep learning  
509 models (reviewed in refs. 188 and 189). These techniques can be model-agnostic or model-  
510 specific, depending on their applicability to all or specific machine learning algorithms, providing  
511 local or global interpretation. Local interpretation focuses on understanding why an ANN gives a  
512 specific output for a given input, whereas global interpretation aims to understand how an ANN  
513 makes predictions using features and parameters learnt from data<sup>190, 191</sup>. LIME (Local  
514 Interpretable Model-agnostic Explanations) is a common local model-agnostic interpretation  
515 method that can explain predictions from ANNs by approximation from interpretable models such  
516 as linear models and decision trees<sup>192</sup>. Other common local interpretation methods include SHAP  
517 (SHapley Additive exPlanations)<sup>193</sup>, based on game theory, and saliency maps<sup>194</sup>, especially for  
518 CNNs. These can provide global interpretation for how neutral and non-neutral variants affect  
519 CNN predictions by aggregating local interpretations across different inputs<sup>75, 79</sup>. Moreover, prior  
520 human knowledge can be useful. For example, deep learning models can be interpreted by finding  
521 correlation between values from hidden units of ANNs and summary statistics from population  
522 genetics theory<sup>106</sup>. Some machine learning algorithms can be explained using population genetics  
523 theory, such as principal component analysis<sup>195, 196</sup>. Therefore, it is interesting to explore ANNs  
524 with population genetics theory and expand population genetics theory with ANNs in the future<sup>197</sup>.

525

## 526 **Conclusions**

527 Although deep learning still has many issues, its power to manage large-scale and multi-modal  
528 data in a multitasking environment<sup>198, 199</sup> is paving a new road to study multiple evolutionary  
529 problems with massive multi-population and multi-omic datasets in this genomic era.  
530 Simultaneously, we believe that biological studies will continue to inspire and improve machine  
531 learning algorithms, just as geneticists contributed to the establishment of linear regression — a  
532 fundamental machine learning algorithm<sup>123</sup> — a century ago<sup>200, 201</sup>.



## 533 References

- 534 1. Wakeley J. The limits of theoretical population genetics. *Genetics* **169**, 1–7 (2005).
- 535 2. Lewontin R. C. Population genetics. *Annu. Rev. Genet.* **1**, 37–70 (1967).
- 536 3. Fu Y. Variances and covariances of linear summary statistics of segregating sites. *Theor.*  
537 *Popul. Biol.* **145**, 95–108 (2022).
- 538 4. Bradburd G. S. & Ralph P. L. Spatial population genetics: It's about time. *Annu. Rev. Ecol.*  
539 *Evol. Syst.* **50**, 427–449 (2019).
- 540 5. Ewens W. J. In *Mathematical Population Genetics I: Theoretical Introduction Second*  
541 *Edition* (Springer, 2004).
- 542 **This classic textbook covers theoretical population genetics ranging from the**  
543 **diffusion theory to the coalescent theory.**
- 544 6. Crow J. F. & Kimura M. In *An Introduction to Population Genetics Theory* (The Blackburn  
545 Press, 2009).
- 546 **This classic textbook introduces the fundamentals of theoretical population**  
547 **genetics.**
- 548 7. Pool J. E., Hellmann I., Jensen J. D. & Nielsen R. Population genetic inference from  
549 genomic sequence variation. *Genome Res.* **20**, 291–300 (2010).
- 550 8. Charlesworth B. & Charlesworth D. Population genetics from 1966 to 2016. *Heredity* **118**,  
551 2–9 (2017).
- 552 9. Johri P. et al. Recommendations for improving statistical inference in population genomics.  
553 *PLoS Biol.* **20**, e3001669 (2022).
- 554 10. The 1000 Genomes Project Consortium. A global reference for human genetic variation.  
555 *Nature* **526**, 68–74 (2015).
- 556 11. Mallick S. et al. The Allen Ancient DNA Resource (AADR): A curated compendium of  
557 ancient human genomes. Preprint at bioRxiv <https://doi.org/10.1101/2023.04.06.535797>  
558 (2023).
- 559 12. The 1001 Genomes Consortium. 1,135 genomes reveal the global pattern of  
560 polymorphism in *Arabidopsis thaliana*. *Cell* **166**, 481–491 (2016).
- 561 13. Sudlow C. et al. UK Biobank: An open access resource for identifying the causes of a wide  
562 range of complex diseases of middle and old Age. *PLoS Med.* **12**, e1001779 (2015).
- 563 14. Walters R. G. et al. Genotyping and population structure of the China Kadoorie Biobank.  
564 Preprint at medRxiv <https://doi.org/10.1101/2022.05.02.22274487> (2022).

- 565 15. Schrider D. R. & Kern A. D. Supervised machine learning for population genetics: A new  
566 paradigm. *Trends Genet.* **34**, 301–312 (2018).
- 567 **This review covers the applications of supervised learning in population genetic**  
568 **inference.**
- 569 16. LeCun Y., Bengio Y. & Hinton G. Deep learning. *Nature* **521**, 436–444 (2015).
- 570 17. Gao H. et al. The landscape of tolerated genetic variation in humans and primates.  
571 *Science* **380**, eabn8153 (2023).
- 572 18. van Hilten A. et al. GenNet framework: Interpretable deep learning for predicting  
573 phenotypes from genetic data. *Commun. Biol.* **4**, 1094 (2021).
- 574 19. Vaswani A. et al. Attention is all you need. In *Advances in Neural Information Processing*  
575 *Systems* (eds. Guyon I. et al.) **30**, 6000–6010 (NIPS, 2017).
- 576 **This study proposed the vanilla transformer architecture, which has become the**  
577 **basis of novel architectures that achieved state-of-the-art performance in different**  
578 **machine learning tasks.**
- 579 20. Sohl-Dickstein J., Weiss E., Maheswaranathan N. & Ganguli S. Deep unsupervised  
580 learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International*  
581 *Conference on Machine Learning* **37**, 2256–2265 (PMLR, 2015).
- 582 21. Nei M. In *Molecular Evolutionary Genetics* 327–403 (Columbia University Press, 1987).
- 583 22. Hamilton M. B. In *Population Genetics* 53–67 (Wiley-Blackwell, 2009).
- 584 23. Kimura M. Diffusion models in population genetics. *J. Appl. Probab.* **1**, 177–232 (1964).
- 585 24. Kingman J. F. C. On the genealogy of large populations. *J. Appl. Probab.* **19**, 27–43  
586 (1982).
- 587 25. Rosenberg N. A. & Nordborg M. Genealogical trees, coalescent theory and the analysis  
588 of genetic polymorphisms. *Nat. Rev. Genet.* **3**, 380–390 (2002).
- 589 26. Fu Y. & Li W. Maximum likelihood estimation of population parameters. *Genetics* **134**,  
590 1261–1270 (1993).
- 591 27. Griffiths R. C. & Tavaré S. Monte Carlo inference methods in population genetics. *Math.*  
592 *Comput. Model.* **23**, 141–158 (1996).
- 593 28. Tavaré S., Balding D. J., Griffiths R. C. & Donnelly P. Inferring coalescence times from  
594 DNA sequence data. *Genetics* **145**, 505–518 (1997).
- 595 29. Marjoram P. & Tavaré S. Modern computational approaches for analysing molecular  
596 genetic variation data. *Nat. Rev. Genet.* **7**, 759–770 (2006).

- 597 30. Williamson S. H. et al. Simultaneous inference of selection and population growth from  
598 patterns of variation in the human genome. *Proc. Natl. Acad. Sci. USA* **102**, 7882–7887  
599 (2005).
- 600 31. Wang M. et al. Detecting recent positive selection with high accuracy and reliability by  
601 conditional coalescent tree. *Mol. Biol. Evol.* **31**, 3068–3080 (2014).
- 602 32. Szpiech Z. A. & Hernandez R. D. selscan: An efficient multithreaded program to perform  
603 EHH-based scans for positive selection. *Mol. Biol. Evol.* **31**, 2824–2827 (2014).
- 604 33. Maclean C. A., Hong N. P. C. & Prendergast J. G. D. hapbin: An efficient program for  
605 performing haplotype-based scans for positive selection in large genomic datasets. *Mol.*  
606 *Biol. Evol.* **32**, 3027–3029 (2015).
- 607 34. Huang X., Kruisz P. & Kuhlwilm M. sstar: A Python package for detecting archaic  
608 introgression from population genetic data with  $S^*$ . *Mol. Biol. Evol.* **39**, msac212 (2022).
- 609 35. Borowiec M. L. et al. Deep learning as a tool for ecology and evolution. *Methods Ecol.*  
610 *Evol.* **13**, 1640–1660 (2022).
- 611 36. Korfmann K., Gaggiotti O. E. & Fumagalli M. Deep learning in population genetics.  
612 *Genome Biol. Evol.* **15**, evad008 (2023).
- 613 37. Alpaydin E. In *Introduction to Machine Learning Third Edition* (The MIT Press, 2014).
- 614 38. Bengio Y., LeCun Y. & Hinton G. Deep learning for AI. *Commun. ACM* **64**, 58–65 (2021).
- 615 39. Sapoval N. et al. Current progress and open challenges for applying deep learning across  
616 the biosciences. *Nat. Commun.* **13**, 1728 (2022).
- 617 40. Bishop C. M. Model-based machine learning. *Phil. Trans. R. Soc. A* **371**, 20120222 (2013).
- 618 41. Lee C., Abdool A. & Huang C. PCA-based population structure inference with generic  
619 clustering algorithms. *BMC Bioinform.* **10**, S73 (2009).
- 620 42. Li H. & Durbin R. Inference of human population history from individual whole-genome  
621 sequences. *Nature* **475**, 493–496 (2011).
- 622 43. Skov L. et al. Detecting archaic introgression using an unadmixed outgroup. *PLoS Genet.*  
623 **14**, e1007641 (2018).
- 624 44. Chen H., Hey J. & Slatkin M. A hidden Markov model for investigating recent positive  
625 selection through haplotype structure. *Theor. Popul. Biol.* **99**, 18–30 (2015).
- 626 45. Lin K., Li H., Schlötterer C. & Futschik A. Distinguishing positive selection from neutral  
627 evolution: Boosting the performance of summary statistics. *Genetics* **187**, 229–244 (2011).
- 628 46. Schrider D. R., Ayroles J., Matute D. R. & Kern A. D. Supervised machine learning reveals  
629 introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLoS Genet.*  
630 **14**, e1007341 (2018).

- 631 47. Durvasula A. & Sankararaman S. A statistical model for reference-free inference of archaic  
632 local ancestry. *PLoS Genet.* **15**, e1008175 (2019).
- 633 48. Goodfellow I., Bengio Y. & Courville A. In *Deep Learning* (MIT Press, 2016).  
634 **This classic textbook introduces the fundamentals of deep learning.**
- 635 49. Eraslan G., Avsec Z., Gagneur J. & Theis F. J. Deep learning: New computational  
636 modelling techniques for genomics. *Nat. Rev. Genet.* **20**, 389–402 (2019).
- 637 50. Villanea F. A. & Schraiber J. G. Multiple episodes of interbreeding between Neanderthals  
638 and modern humans. *Nat. Ecol. Evol.* **3**, 39–44 (2019).
- 639 51. Unadkat S. B., Ciocoiu M. M. & Medsker L. R. Introduction. In *Recurrent neural networks:  
640 Design and applications* (eds. Medsker L. R. & Jain L. C.) 1–12 (CRC Press, 1999).
- 641 52. Géron A. Introduction to artificial neural networks. In *Neural networks and deep learning*.  
642 <https://www.oreilly.com/library/view/neural-networks-and/9781492037354/ch01.html>  
643 (O'Reilly Media, Inc. 2018).
- 644 53. Sheehan S. & Song Y. S. Deep learning for population genetic inference. *PLoS Comput.  
645 Biol.* **12**, e1004845 (2016).
- 646 54. Mondal M., Bertranpetit J. & Lao O. Approximate Bayesian computation with deep learning  
647 supports a third archaic introgression in Asia and Oceania. *Nat. Commun.* **10**, 246 (2019).
- 648 55. Sanchez T, Curry J. & Jay F. Deep learning for population size history inference: Design,  
649 comparison and combination with approximate Bayesian computation. *Mol. Ecol. Resour.*  
650 **21**, 2645–2660 (2021).
- 651 56. Tran L. N., Sun C. K., Struck T. J., Sajan M. & Gutenkunst R. N. Computationally efficient  
652 demographic history inference from allele frequencies with supervised machine learning.  
653 Preprint at bioRxiv <https://doi.org/10.1101/2023.05.24.542158> (2023).
- 654 57. Romero A. et al. Diet networks: Thin parameters for fat genomics. In *5th International  
655 Conference on Learning Representations (ICLR, 2017)*.
- 656 58. Isildak U., Stella A. & Fumagalli M. Distinguishing between recent balancing selection and  
657 incomplete sweep using deep neural networks. *Mol. Ecol. Resour.* **21**, 2706–2718 (2021).
- 658 59. Qin X., Chiang C. W. K. & Gaggiotti O. E. Deciphering signatures of natural selection via  
659 deep learning. *Brief. Bioinformatics* **23**, 1–10 (2022).
- 660 60. Burger K. E., Pfaffelhuber P. & Baumdicker F. Neural networks for self-adjusting mutation  
661 rate estimation when the recombination rate is unknown. *PLoS Comput. Biol.* **18**,  
662 e1010407 (2022).
- 663 61. Fang Y., Deng S. & Li C. A generalizable deep learning framework for inferring fine-scale  
664 germline mutation rate maps. *Nat. Mach. Intell.* **4**, 1209–1223 (2022).

- 665 62. Battey C. J., Ralph P. L. & Kern A. D. Predicting geographic location from genetic variation  
666 with deep neural networks. *eLife* **9**, e54507 (2020).
- 667 63. Flagel L., Brandvain Y. & Schrider D. R. The unreasonable effectiveness of convolutional  
668 neural networks in population genetic inference. *Mol. Biol. Evol.* **36**, 220–238 (2019).  
669 **This study experimented convolutional neural networks for various tasks in**  
670 **population genetic inference.**
- 671 64. Wang Z. et al. Automatic inference of demographic parameters using generative  
672 adversarial network. *Mol. Ecol. Resour.* **21**, 2689–2705 (2021).  
673 **This study developed a generative adversarial framework aimed at inferring**  
674 **demographic parameters from data in an unsupervised manner.**
- 675 65. Silver D. et al. Mastering the game of Go with deep neural networks and tree search.  
676 *Nature* **529**, 484–489 (2016).
- 677 66. Montserrat D. M., Bustamante C. & Ioannidis A. LAI-Net: Local-ancestry inference with  
678 neural networks. In *International Conference on Acoustics, Speech, & Signal Processing*  
679 *2020*, 1314–1318 (ICASSP, 2020).
- 680 67. Sabat B. O., Montserrat D. M., Giró-i-Nieto X. & Ioannidis A. G. SALAI-Net: Species-  
681 agnostic local ancestry inference network. *Bioinformatics* **38**, ii27–ii33 (2022).
- 682 68. Kern A. D. & Schrider D. R. diploS/HIC: An updated approach to classifying selective  
683 sweeps. *G3*, **8**: 1959–1970 (2018).
- 684 69. Torada L. et al. ImaGene: A convolutional neural network to quantify natural selection from  
685 population genomic data. *BMC Bioinform.* **20**, 337 (2019).
- 686 70. Deelder W. et al. Using deep learning to identify recent positive selection in malaria  
687 parasite sequence data. *Malar. J.* **20**, 270 (2021).
- 688 71. Xue A. T., Schrider D. R., Kern A. D. & Ag1000g Consortium. Discovery of ongoing  
689 selective sweeps within Anopheles mosquito populations using deep learning. *Mol. Biol.*  
690 *Evol.* **38**, 1168–1183 (2021).
- 691 72. Caldas I. V., Clark A. G. & Messer P. W. Inference of selective sweep parameters through  
692 supervised learning. Preprint at bioRxiv <https://doi.org/10.1101/2022.07.19.500702>  
693 (2022).
- 694 73. Hamid I, Korunes K. L., Schrider D. R. & Goldberg A. Localizing post-admixture adaptive  
695 variants with object detection on ancestry-painted chromosomes. *Mol. Biol. Evol.* **40**,  
696 msad074 (2023).
- 697 74. Whitehouse L. S. & Schrider D. R. Timesweeper: Accurately identifying selective sweeps  
698 using population genomic time series. *Genetics* **224**, iyad084 (2023).

- 699 75. Cecil R. M. & Sugden L. A. On convolutional neural networks for selection inference:  
700 Revealing the lurking role of preprocessing, and the surprising effectiveness of summary  
701 statistics. Preprint at bioRxiv <https://doi.org/10.1101/2023.02.26.530156> (2023).
- 702 76. Arnab S. P., Amin M. R. & DeGiorgio M. Uncovering footprints of natural selection through  
703 time-frequency analysis of genomic summary statistics. Preprint at bioRxiv  
704 <https://doi.org/10.1101/2022.10.05.510997> (2022).
- 705 77. Lauterbur M. E., Munch K. & Enard D. Versatile detection of diverse selective sweeps with  
706 Flex-sweep. *Mol. Biol. Evol.* **40**, msad139 (2023).
- 707 78. Blischak P. D., Barker M. S. & Gutenkunst R. N. Chromosome-scale inference of hybrid  
708 speciation and admixture with convolution neural networks. *Mol. Ecol. Resour.* **21**, 2676–  
709 2688 (2021).
- 710 79. Gower G., Picazo P. I., Fumagalli M. & Racimo F. Detecting adaptive introgression in  
711 human evolution using convolutional neural networks. *eLife* **10**, e64669 (2021).
- 712 80. Ray D. D., Flagel L., & Schrider D. R. IntroUNET: Identifying introgressed alleles via  
713 semantic. Preprint at bioRxiv <https://doi.org/10.1101/2023.02.07.527435> (2023).
- 714 81. Zhang Y. et al. Inferring historical introgression with deep learning. *Syst. Biol.*  
715 <https://doi.org/10.1093/sysbio/syad033> (2023).
- 716 82. Smith C. C. R., Tittes S., Ralph P. L. & Kern A. D. Dispersal inference from population  
717 genetic variation using a convolutional neural network. *Genetics* **224**, iyad068 (2023).
- 718 83. Battey C. J., Coffing G. C. & Kern A. D. Visualizing population structure with variational  
719 autoencoders. *G3* **11**, jkaa036 (2021).
- 720 84. Ausmees K. & Nettelblad C. A deep learning framework for characterization of genotype  
721 data. *G3* **12**, jkac020 (2022).
- 722 85. Booker W. W., Ray D. D. & Schrider D. R. This population doesn't exist: Learning the  
723 distribution of evolutionary histories with generative adversarial networks. *Genetics* **224**,  
724 iyad063 (2023).
- 725 86. Meisner J. & Albrechtsen A. Haplotype and population structure inference using neural  
726 networks in whole-genome sequencing data. *Genome Res.* **32**, 1542–1552 (2022).
- 727 **This study developed a variational autoencoder scalable on the UK Biobank dataset**  
728 **for estimating ancestry proportions across the genome without training from**  
729 **simulated data.**
- 730 87. Yelmen B. et al. Deep convolutional and conditional neural networks for large-scale  
731 genomic data generation. Preprint at bioRxiv <https://doi.org/10.1101/2023.03.07.530442>  
732 (2023).

- 733 88. Chan J. et al. A likelihood-free inference framework for population genetic data using  
734 exchangeable neural networks. In *Advances in Neural Information Processing Systems*  
735 (eds. Bengio S. et al.) **31**, 8594–8605 (NeurIPS, 2018).
- 736 89. Graves A. In *Supervised Sequence Labelling with Recurrent Neural Networks* (Springer  
737 2012).
- 738 90. Adrion J. R., Galloway J. G. & Kern A. D. Predicting the landscape of recombination using  
739 deep learning. *Mol. Biol. Evol.* **37**, 1790–1808 (2020).
- 740 91. Hejase H. A., Mo Z., Campagna L. & Siepel A. A deep-learning approach for inference of  
741 selective sweeps from ancestral recombination graph. *Mol. Biol. Evol.* **39**, msab332  
742 (2022).
- 743 92. Sanchez-Lengeling B., Reif E., Pearce A. & Wiltchko A. B. A gentle introduction to graph  
744 neural networks. *Distill* <https://doi.org/10.23915/distill.00033> (2021).
- 745 93. Daigavane A., Ravindran B. & Aggarwal G. Understanding convolutions on graphs. *Distill*  
746 <https://doi.org/10.23915/distill.00032> (2021).
- 747 94. Veličković P. et al. Graph attention networks. In *6th International Conference on Learning*  
748 *Representations* (ICLR, 2018).
- 749 95. Griffiths R. C. & Marjoram P. Ancestral inference from samples of DNA sequences with  
750 recombination. *J. Comput. Biol.* **3**, 479–502 (1996).
- 751 96. Paradis E. Analysis of haplotype networks: The randomized minimum spanning tree  
752 method. *Methods Ecol. Evol.* **9**, 1308–1317 (2018).
- 753 97. Korfmann K., Sellinger T., Freund F., Fumagalli M. & Tellier A. Simultaneous inference of  
754 past demography and selection from the ancestral recombination graph under the beta  
755 coalescent. Preprint at bioRxiv <https://doi.org/10.1101/2022.09.28.508873> (2022).
- 756 98. Hinton G. et al. Deep neural networks for acoustic modeling in speech recognition: The  
757 shared views of four research groups. In *IEEE Signal Process. Mag.* **29**, 82–97 (2012).
- 758 99. Bond-Taylor S., Leach A., Long Y. & Willcocks C. G. Deep generative modelling: A  
759 comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive  
760 models. *IEEE Trans. Pattern Anal. Mach. Intell.* **44**, 7327–7346 (2022).
- 761 100. Yelmen B. et al. Creating artificial human genomes using generative neural networks.  
762 *PLoS Genet.* **17**, e1009303 (2021).
- 763 **This study utilized restricted Boltzmann machines and generative adversarial**  
764 **networks for synthesizing realistic human genomes.**

- 765 101. Goodfellow I. J. et al. Generative adversarial nets. In *Advances in Neural Information*  
766 *Processing Systems* (eds. Ghahramani Z., Welling M., Cortes C., Lawrence N. &  
767 Weinberger K. Q.) **27**, 2672–2680 (NIPS, 2014).
- 768 102. Saxena D. & Cao J. Generative adversarial networks (GANs): Challenges, solutions, and  
769 future directions. *ACM Comput. Surv.* **54**, 1–42 (2021).
- 770 103. Mantes A. D., Montserrat D. M., Bustamante C. D., Giró-i-Nieto X. & Ioannidis A. G. Neural  
771 ADMIXTURE: Rapid population clustering with autoencoders. *Nat. Comput. Sci.*  
772 <https://doi.org/10.1038/s43588-023-00482-7> (2023).
- 773 104. Alexander D. H., Novembre J. & Lange K. Fast model-based estimation of ancestry in  
774 unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 775 105. Lawson D. J., Hellenthal G., Myers S. & Falush D. Inference of population structure using  
776 dense haplotype data. *PLoS Genet.* **8**, e1002453 (2012).
- 777 106. Riley R., Mathieson I. & Mathieson S. Interpreting generative adversarial networks to infer  
778 natural selection from genetic data. Preprint at bioRxiv  
779 <https://doi.org/10.1101/2023.03.07.531546> (2023).
- 780 107. Gower G., Picazo P. I., Lindgren F. & Racimo F. Inference of population genetics  
781 parameters using discriminator neural networks: An adversarial Monte Carlo approach.  
782 Preprint at bioRxiv <https://doi.org/10.1101/2023.04.27.538386> (2023).
- 783 108. Montserrat D. M., Bustamante C. & Ioannidis A. Class-conditional VAE-GAN for local-  
784 ancestry simulation. Preprint at arXiv <https://doi.org/10.48550/arXiv.1911.13220> (2019).
- 785 109. Borji A. Pros and cons of GAN evaluation measures. *Comput. Vis. Image Underst.* **179**,  
786 41–65 (2019).
- 787 110. Phuong M. & Hutter M. Formal algorithms for transformers. Preprint at arXiv  
788 <https://arxiv.org/abs/2207.09238> (2022).
- 789 111. Katharopoulos A., Vyas A., Pappas N. & Fleuret F. Transformers are RNNs: Fast  
790 autoregressive transformers with linear attention. In *Proceedings of the 37th International*  
791 *Conference on Machine Learning* **119**, 5156–5165 (PMLR, 2020).
- 792 112. Cordonnier J., Loukas A. & Jaggi M. On the relationship between self-attention and  
793 convolutional layers. Preprint at arXiv <https://doi.org/10.48550/arXiv.1911.03584> (2019).
- 794 113. Lakew S. M., Cettolo M. & Federico M. A comparison of transformer and recurrent neural  
795 networks on multilingual neural machine translation. In *Proceedings of the 27th*  
796 *International Conference on Computational Linguistics (ICCL)*, (2018).
- 797 114. Ramachandran P. et al. Stand-alone self-attention in vision models. In *Advances in Neural*  
798 *Information Processing Systems* (eds. Wallach H. et al.) **32**, 68–80 (NeurIPS, 2019).



- 799 115. Liu Y. X. et al. Learning virus genotype-fitness landscape in embedding space. Preprint at  
800 bioRxiv <https://doi.org/10.1101/2023.02.09.527693> (2023).
- 801 116. Devlin J., Chang M., Lee K. & Toutanova K. BERT: Pre-training of deep bidirectional  
802 transformers for language understanding. In *Proceedings of the 2019 Conference of the*  
803 *North American Chapter of the Association for Computational Linguistics: Human*  
804 *Language Technologies 1*, 4171–4186 (ACL, 2019).
- 805 117. Brown T. B. et al. Language models are few-shot learners. In *Advances in Neural*  
806 *Information Processing Systems* (eds. Larochelle H., Ranzato M., Hadsell R., Balcan M.  
807 F. & Lin H.-T.) **33**, 1877–1901 (NeurIPS, 2020).
- 808 118. Dosovitskiy A. et al. An image is worth 16x16 words: Transformers for image recognition  
809 at scale. In *9th International Conference on Learning Representations (ICLR, 2021)*.
- 810 119. Zaheer M. et al. Big bird: Transformers for longer sequences. In *Advances in Neural*  
811 *Information Processing Systems* (eds. Larochelle H., Ranzato M., Hadsell R., Balcan M.  
812 F. & Lin H.-T.) **33**, 17283–17297 (NeurIPS, 2020).
- 813 120. Dhariwal P. & Nichol A. Q. Diffusion models beat GANs on image synthesis. In *Advances*  
814 *in Neural Information Processing Systems* (eds. Ranzato M., Beygelzimer A., Dauphin Y.  
815 N., Liang P. S. & Vaughan J. W.) **34**, 8780–8794 (NeurIPS, 2021).
- 816 121. Croitoru F., Hondru V. & Ionescu R. T. Diffusion models in vision: A survey. Preprint at  
817 arXiv <https://doi.org/10.48550/arXiv.2209.04747> (2022).
- 818 122. Huang Y. & Siepel A. Estimation of allele-specific fitness effects across human protein-  
819 coding sequences and implications for disease. *Genome Res.* **29**, 1310–1321 (2019).
- 820 123. Bishop C. M. In *Pattern Recognition and Machine Learning* (Springer, 2006).
- 821 **This classic textbook covers a range of machine learning algorithms and statistical**  
822 **inference approaches, which are also widely used in population genetic inference.**
- 823 124. Bengio Y. Practical recommendations for gradient-based training of deep architectures. In  
824 *Neural Networks: Tricks of the Trade* (eds. Montavon G., Orr G. B. & Müller K. R.), 437–  
825 478 (Springer, 2012).
- 826 125. Tieleman T. & Hinton G. Lecture 6.5-rmsprop: Divide the gradient by a running average of  
827 its recent magnitude. In *Coursera: Neural networks for machine learning 4*, 26–31  
828 (Coursera, 2012).
- 829 126. Kingma D. & Ba J. Adam: A method for stochastic optimization. In *3rd International*  
830 *Conference on Learning Representations (ICLR, 2015)*.

- 831 127. Jospin L. V., Laga H., Boussaid F., Buntine W. & Bennamoun M. Hands-on Bayesian  
832 neural networks — A tutorial for deep learning users. *IEEE Comput. Intell. Mag.* **17**, 29–48  
833 (2022).
- 834 128. Srivastava N., Hinton G., Krizhevsky A., Sutskever I. & Salakhutdinov R. Dropout: A simple  
835 way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958  
836 (2014).
- 837 129. Prechelt L. Early stopping — But when? In *Neural Networks: Tricks of the Trade* (eds.  
838 Montavon G., Orr G. B. & Müller K. R.), 53–67 (Springer, 2012).
- 839 130. Arlot S. & Celisse A. A survey of cross-validation procedures for model selection. *Stat.*  
840 *Surv.* **4**, 40–79 (2010).
- 841 131. Ioffe S. & Szegedy C. Batch normalization: Accelerating deep network training by reducing  
842 internal covariate shift. In *Proceedings of the 32nd International Conference on  
843 International Conference on Machine Learning* **37**, 448–456 (PMLR, 2015).
- 844 132. Luo P., Wang X., Shao W. & Peng Z. Towards understanding regularization in batch  
845 normalization. In *7th International Conference on Learning Representations* (ICLR, 2019).
- 846 133. Green R. E. et al. A draft sequence of the Neanderthal genome. *Science* **328**, 710–722  
847 (2010).
- 848 134. Borji A. Pros and cons of GAN evaluation measures: New developments. *Comput. Vis.*  
849 *Image Underst.* **215**, 103329 (2022).
- 850 135. Theis L., van den Oord A. & Bethge M. A note on the evaluation of generative models. In  
851 *4th International Conference on Learning Representations* (ICLR, 2016).
- 852 136. Sajjadi M. S. M., Bachem O., Lucic M., Bousquet O. & Gelly S. Assessing generative  
853 models via precision and recall. In *Advances in Neural Information Processing Systems*  
854 (eds. Bengio S. et al.) **31**, 5234–5243 (NeurIPS, 2018).
- 855 137. Naeem M. F., Oh S. J., Uh Y., Choi Y. & Yoo J. Reliable fidelity and diversity metrics for  
856 generative models. In *Proceedings of the 37th International Conference on Machine  
857 Learning* **119**, 7176–7185 (PMLR, 2020).
- 858 138. Perera M. et al. Generative moment matching networks for genotype simulation. In *44th  
859 Annual International Conference of the IEEE Engineering in Medicine & Biology Society*,  
860 1379–1383 (EMBC, 2022).
- 861 139. Kynkäänniemi T., Karras T., Laine S., Lehtinen J. & Aila T. Improved precision and recall  
862 metric for assessing generative models. In *Advances in Neural Information Processing  
863 Systems* (eds. Wallach H. et al.) **32**, 3927–3936 (NeurIPS, 2019).

- 864 140. Cornuet J. M., Aulagnier S., Lek S., Franck S. & Solignac M. Classifying individuals among  
865 infra-specific taxa using microsatellite data and neural networks. *C. R. Acad. Sci. III.* **319**,  
866 1167–1177 (1996).
- 867 141. Guinand B., Topchy A., Page K. S., Burnham-Curtis M. K., Punch W. F. & Scribner K. T.  
868 Comparisons of likelihood and machine learning methods of individual classification. *J.*  
869 *Hered.* **93**, 260–269 (2002).
- 870 142. Sengupta S. et al. A review of deep learning with special emphasis on architectures,  
871 applications and recent trends. *Knowl. Based Syst.* **194**, 105596 (2020).
- 872 143. Hornik K., Stinchcombe M. & White H. Multilayer feedforward networks are universal  
873 approximators. *Neural Netw.* **2**, 359–366 (1989).
- 874 144. Schäfer A. M. & Zimmermann H. G. Recurrent neural networks are universal  
875 approximators. In *Artificial Neural Networks – ICANN 2006* (eds. Kollias S. D., Stafylopatis  
876 A., Duch W. & Oja E.), 632–640 (Springer, 2006).
- 877 145. Browning S. R., Browning B. L., Zhou Y., Tucci S. & Akey J. M. Analysis of human  
878 sequence data reveals two pulses of archaic Denisovan admixture. *Cell* **173**, 53–61  
879 (2018).
- 880 146. Frolov S., Hinz T., Raue F., Hees J. & Dengel A. Adversarial text-to-image synthesis: A  
881 review. *Neural Netw.* **144**, 187–209 (2021).
- 882 147. Abrantes J. P., Abrantes A. J. & Oliehoek F. A. Mimicking evolution with reinforcement  
883 learning. Preprint at arXiv <https://arxiv.org/abs/2004.00048> (2020).
- 884 148. Fawzi A. et al. Discovering faster matrix multiplication algorithms with reinforcement  
885 learning. *Nature* **610**, 47–53 (2022).
- 886 149. Mankowitz D. J. et al. Faster sorting algorithms discovered using deep reinforcement  
887 learning. *Nature* **618**, 257–263 (2023).
- 888 150. Hui Z., Li J., Wang X. & Gao X. Learning the non-differentiable optimization for blind super-  
889 resolution. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2093–  
890 2102 (CVPR, 2021).
- 891 151. Silver D. et al. Mastering the game of Go without human knowledge. *Nature* **550**, 354–359  
892 (2017).
- 893 152. Ibnu C. R. M., Santoso J. & Surendro K. Determining the neural Network topology: A  
894 Review. In *Proceedings of the 8th International Conference on Software and Computer*  
895 *Applications*, 357–362 (ICSCA, 2019).
- 896 153. Menghani G. Efficient deep learning: A survey on making deep learning models smaller,  
897 faster, and better. *ACM Comput. Surv.* **55**, 1–37 (2023).

- 898 154. He K., Zhang X., Ren S. & Sun J. Identity mappings in deep residual networks. In  
899 *Computer Vision – ECCV 2016* (eds. Leibe B., Matas J. Sebe N. & Welling M.) 630–645  
900 (Springer, 2016).
- 901 155. Brown T. B. et al. Language models are few-shot learners. In *Advances in Neural*  
902 *Information Processing Systems* (eds. Larochelle H., Ranzato M., Hadsell R., Balcan M.  
903 F. & Lin H.-T.) **33**, 1877–1901 (NeurIPS, 2020).
- 904 156. Ouyang L. et al. Training language models to follow instructions with human feedback.  
905 Preprint at arXiv <https://doi.org/10.48550/arXiv.2203.02155> (2022).
- 906 157. Touvron H. et al. LLaMA: Open and efficient foundation language models. Preprint at arXiv  
907 <https://doi.org/10.48550/arXiv.2302.13971> (2023).
- 908 158. Kang M. et al. Scaling up GANs for text-to-image synthesis. Preprint at arXiv  
909 <https://doi.org/10.48550/arXiv.2303.05511> (2023).
- 910 159. Kao W.-T. & Lee H.-Y. Is BERT a cross-disciplinary knowledge learner? A surprising  
911 finding of pre-trained models' transferability. In *Findings of the Association for*  
912 *Computational Linguistics: EMNLP 2021*, 2195–2208 (ACL, 2021).
- 913 160. Marinó G. C., Petrini A., Malchiodi D. & Frasca M. Deep neural networks compression: A  
914 comparative survey and choice recommendations. *Neurocomputing* **520**, 152–170 (2023).
- 915 161. Finn C., Abbeel P. & Levine S. Model-agnostic meta-learning for fast adaptation of deep  
916 networks. In *Proceedings of the 34th International Conference on Machine Learning* (eds.  
917 Precup D. & Teh Y. W.) **70**, 1126–1135 (PMLR, 2017).
- 918 162. Wei Y., Zhao P. & Huang J. Meta-learning hyperparameter performance prediction with  
919 neural processes. In *Proceedings of the 34th International Conference on Machine*  
920 *Learning* (eds. Meila M. & Zhang T.) **139**, 11058–11067 (PMLR, 2021).
- 921 163. Hospedales T., Antoniou A., Micaelli P. & Storkey A. Meta-learning in neural networks: A  
922 survey. *IEEE Trans. Pattern Anal. Mach. Intell.* **44**, 5149–5169 (2022).
- 923 164. Kaveh M. & Mesgari M. S. Application of meta-heuristic algorithms for training neural  
924 networks and deep learning architectures: A comprehensive review. *Neural Process. Lett.*  
925 <https://doi.org/10.1007/s11063-022-11055-6> (2022).
- 926 165. Tirumala S. S., Ali S. & Ramesh C. P. Evolving deep neural networks: A new prospect. In  
927 *12th International Conference on Natural Computation, Fuzzy Systems and Knowledge*  
928 *Discovery*, 69–74 (ICNC-FSKD, 2016).
- 929 166. Stanley K. O., Clune J., Lehman J. & Miikkulainen R. Designing neural networks through  
930 neuroevolution. *Nat. Mach. Intell.* **1**, 24–35 (2019).

- 931 167. Juan D., Santpere G., Kelley J. L., Cornejo O. E. & Marques-Bonet T. Current advances  
932 in primate genomics: Novel approaches for understanding evolution and disease. *Nat.*  
933 *Rev. Genet.* **24**, 314–331 (2023).
- 934 168. Wang Y., Yao Q., Kwok J. T. & Ni L. M. Generalizing from a few examples: A survey on  
935 few-shot learning. *ACM Comput. Surv.* **53**, 1–34 (2020).
- 936 169. Wang W., Zheng V. W., Yu H. & Miao C. A survey of zero-shot learning: Settings, methods,  
937 and applications. *ACM Trans. Intell. Syst. Technol.* **10**, 1–37 (2019).
- 938 170. Saada J. N., Hu A. & Palamara P. F. Inference of pairwise coalescence times and allele  
939 ages using deep neural networks. In *Workshop on Learning Meaningful Representations*  
940 *of Life at the 35th Conference on Neural Information Processing Systems*  
941 <https://www.lmrl.org/papers2021> (LMRL, 2021).
- 942 171. Lauterbur M. E. et al. Expanding the stdpopsim species catalog, and lessons learned for  
943 realistic genome simulations. *eLife* **12**, RP84874 (2023).
- 944 172. Hudson R. R. Generating samples under a Wright-Fisher neutral model of genetic  
945 variation. *Bioinformatics* **18**, 337–338 (2002).
- 946 173. Baumdicker F. et al. Efficient ancestry and mutation simulation with msprime 1.0. *Genetics*  
947 **220**, iyab229 (2022).
- 948 174. Haller B. C. & Messer P. W. SLiM 3: Forward genetic simulations beyond the Wright-Fisher  
949 model. *Mol. Biol. Evol.* **36**, 632–637 (2019).
- 950 175. Huang X. et al. Inferring genome-wide correlation of mutation fitness effects between  
951 populations. *Mol. Biol. Evol.* **38**, 4588–4602 (2021).
- 952 176. Ewing G. B. & Jensen J. D. The consequences of not accounting for background selection  
953 in demographic inference. *Mol. Ecol.* **25**, 135–141 (2016).
- 954 177. Mo Z. & Siepel A. Domain-adaptive neural networks improve supervised machine learning  
955 based on simulated population genetic data. Preprint at bioRxiv  
956 <https://doi.org/10.1101/2023.03.01.529396> (2023).
- 957 178. Hendrycks D., Lee K. & Mazeika M. Using pre-training can improve model robustness and  
958 uncertainty. In *Proceedings of the 36th International Conference on Machine Learning* **97**,  
959 2712–2721 (PMLR, 2019).
- 960 179. Hendrycks D., Mazeika M., Kadavath S. & Song D. Using self-supervised learning can  
961 improve model robustness and uncertainty. In *Advances in Neural Information Processing*  
962 *Systems* (eds. Wallach H. et al.) **32**, 15663–15674 (NeurIPS, 2019).
- 963 180. Abadi M. et al. TensorFlow: Large-scale machine learning on heterogeneous systems.  
964 *TensorFlow*.

- 965 <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/45166>  
966 [.pdf](#) (2015).
- 967 181. Paszke A. et al. PyTorch: An imperative style, high-performance deep learning library. In  
968 *Advances in Neural Information Processing Systems* (eds. Wallach H. M. et al.) **32**, 8026–  
969 8037 (NeurIPS, 2019).
- 970 182. Chen B. et al. Towards training reproducible deep learning models. In *Proceedings of the*  
971 *44th International Conference on Software Engineering*, 2202–2214 (ACM, 2022).
- 972 183. Walsh I. et al. DOME: Recommendations for supervised machine learning validation in  
973 biology. *Nat. Methods* **18**, 1122–1127 (2021).
- 974 184. Sanchez T. et al. dnadna: A deep learning framework for population genetics inference.  
975 *Bioinformatics* **39**, btac765 (2023).
- 976 185. Montserrat D. M. & Ioannidis A. G. Adversarial attacks on genotype sequences. Preprint  
977 at bioRxiv <https://doi.org/10.1101/2022.11.07.515527> (2022).
- 978 186. Ren K., Zheng T., Qin Z. & Liu X. Adversarial attacks and defenses in deep learning.  
979 *Engineering* **6**, 346–360 (2020).
- 980 187. Montavon G., Samek W. & Müller K. Methods for interpreting and understanding deep  
981 neural networks. *Digit. Signal Process.* **73**, 1–15 (2018).
- 982 188. Azodi C. B., Tang J. & Shiu S. Opening the black box: Interpretable machine learning for  
983 geneticists. *Trends Genet.* **36**, 442–455 (2020).
- 984 189. Novakovsky G., Dexter N., Libbrecht M. W., Wasserman W. W. & Mostafavi S. Obtaining  
985 genetics insights from deep learning via explainable artificial intelligence. *Nat. Rev. Genet.*  
986 **24**, 125–137 (2023).
- 987 190. Liang Y., Li S., Yan C., Li M. & Jiang C. Explaining the black-box model: A survey of local  
988 interpretation methods for deep neural networks. *Neurocomputing* **419**, 168–182 (2021).
- 989 191. Saleem R., Yuan B., Kurugollu F., Anjum A. & Liu L. Explaining deep neural networks: A  
990 survey on the global interpretation methods. *Neurocomputing* **513**, 165–180 (2022).
- 991 192. Ribeiro M. T., Singh S. & Guestrin C. “Why should I trust you?” Explaining the predictions  
992 of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on*  
993 *Knowledge Discovery and Data Mining*, 1135–1144 (ACM, 2016).
- 994 193. Lundberg S. M. & Lee S. A unified approach to interpreting model predictions. In *Advances*  
995 *in Neural Information Processing Systems* (eds. Guyon I. et al.) **30**, 4765–4774 (NIPS,  
996 2017).

- 997 194. Simonyan K., Vedaldi A. & Zisserman A. Deep inside convolutional networks: Visualising  
998 image classification models and saliency maps. Preprint at arXiv  
999 <https://doi.org/10.48550/arXiv.1312.6034> (2013).
- 1000 195. McVean G. A genealogical interpretation of principal components analysis. *PLoS Genet.*  
1001 **5**, 1000686 (2009).
- 1002 196. Peter B. M. A geometric relationship of  $F_2$ ,  $F_3$  and  $F_4$ -statistics with principal component  
1003 analysis. *Phil. Trans. R. Soc. B* **377**, 20200413 (2022).
- 1004 197. Tenachi W., Ibata R. & Diakogiannis F. Deep symbolic regression for physics guided by  
1005 units constraints: Toward the automated discovery of physical laws. Preprint at arXiv  
1006 <https://doi.org/10.48550/arXiv.2303.03192> (2023).
- 1007 198. OpenAI. GPT-4 technical report. *OpenAI*. <https://cdn.openai.com/papers/gpt-4.pdf> (2023).
- 1008 199. Bubeck S. et al. Sparks of artificial general intelligence: Early experiments with GPT-4.  
1009 Preprint at arXiv <https://doi.org/10.48550/arXiv.2303.12712> (2023).
- 1010 200. Pearson K. Notes on the history of correlation. *Biometrika* **13**, 25–45 (1920).
- 1011 201. Denis D. J. The origins of correlation and regression: Francis Galton or Auguste Bravais  
1012 and the error theorists? *Hist. Philos. Psychol. Bull.* **13**, 36–44 (2001).
- 1013 202. Ho J., Jain A. & Abbeel P. Denoising diffusion probabilistic models. In *Advances in Neural*  
1014 *Information Processing Systems* (eds. Larochelle H., Ranzato M., Hadsell R., Balcan M.  
1015 F. & Lin H.-T.) **33**, 6840–6851 (NeurIPS, 2020).
- 1016 203. Patel A., Montserrat D. M., Bustamante C. D. & Ioannidis A. Hyperbolic geometry-based  
1017 deep learning methods to produce population trees from genotype data. Preprint at  
1018 bioRxiv <https://doi.org/10.1101/2022.03.28.484797> (2022).
- 1019

## 1020 Acknowledgements

1021 The authors are grateful for the support from the Life Science Compute Cluster (LiSC) of the  
1022 University of Vienna and the feedback from R.N. Gutenkunst. X.H. thanks H.-T. Lin and H.-Y.  
1023 Lee, who provided extraordinary open courses online for learning machine learning. O.D.  
1024 acknowledges support from the John Templeton Foundation (Id: 62178). M.K. is funded by the  
1025 Vienna Science and Technology Fund (WWTF) [10.47379/VRG20001]. ChatGPT with GPT-4  
1026 from OpenAI was utilized for language editing.

1027

## 1028 **Author contributions**

1029 X.H., M.K. and O.L. researched the literature and wrote the article. All authors substantially  
1030 contributed to discussions of the content and reviewed and/or edited the manuscript before  
1031 submission.

## 1032 **Competing interests**

1033 The authors declare no competing interests.

1034

## 1035 **Peer review information**

1036 *Nature Reviews Genetics* thanks Alexander Ioannidis, and the other, anonymous, reviewer(s) for  
1037 their contribution to the peer review of this work.

1038

## 1039 **Related links**

1040 Open Neural Network Exchange: <https://onnx.ai/>



**Table 1. Recent open-source deep learning software for population genetic inference**

Tool	Population genetics problem	Learning paradigm	Architecture	Programming language	Deep learning framework or library	Available pre-trained model
LAI-Net <sup>66</sup>	Admixture	Supervised	CNN	Python	PyTorch	Yes
Neural ADMIXTURE <sup>103</sup>	Admixture	Unsupervised	AE	Python	PyTorch	No
HaploNet <sup>86</sup>	Admixture	Unsupervised	VAE	Python	PyTorch	No
SALAI-Net <sup>67</sup>	Admixture	Supervised	CNN	Python	PyTorch	Yes
donni <sup>56</sup>	Demography	Supervised	FNN	Python	Scikit-learn	Yes
pg-gan <sup>64</sup>	Demography	Unsupervised	GAN-like	Python	TensorFlow	No
evoNet <sup>53</sup>	Demography Natural selection	Supervised	FNN	Java	Customized	No
genomatnn <sup>79</sup>	Introgression	Supervised	CNN	Python	TensorFlow	Yes
introUNET <sup>80</sup>	Introgression	Supervised	CNN	Python	PyTorch	No
ERICA <sup>81</sup>	Introgression	Supervised	CNN	Python	TensorFlow	Yes
MuRaL <sup>61</sup>	Mutation rate	Supervised	FNN/CNN	Python	PyTorch	Yes
diploS/HIC <sup>68</sup>	Natural selection	Supervised	CNN	Python	TensorFlow	No
Flex-sweep <sup>77</sup>	Natural selection	Supervised	CNN	Python	TensorFlow	Yes
Timesweeper <sup>74</sup>	Natural selection	Supervised	CNN	Python	TensorFlow	No
disc-pg-gan <sup>106</sup>	Natural selection	Unsupervised	GAN-like	Python	TensorFlow	No
popvae <sup>83</sup>	Population structure	Unsupervised	VAE	Python	TensorFlow	No
GenoCAE <sup>84</sup>	Population structure	Unsupervised	AE	Python	TensorFlow	No

ReLERNN <sup>90</sup>	Recombination rate	Supervised	RNN	Python	TensorFlow	No
PG-Alignments-GAN <sup>85</sup>	Simulation	Unsupervised	GAN	Python	PyTorch	No
Locator <sup>62</sup>	Spatial pattern	Supervised	FNN	Python	TensorFlow	No
disperseNN <sup>82</sup>	Spatial pattern	Supervised	CNN	Python	TensorFlow	Yes

1042 AE, autoencoder; CNN, convolutional neural network; FNN, feed-forward neural network; GAN,  
1043 generative adversarial network; RNN, recurrent neural network; VAE, variational autoencoder.  
1044 We have categorized open-source software developed since 2016 that employs common artificial  
1045 neural network architectures based on population genetics problems they address. We define a  
1046 tool as software if it includes a detailed manual and command-line interface for user-defined data  
1047 and parameters. We have excluded software for improving deep learning model performance and  
1048 implementation from this table, as they are not specific to any population genetics problems<sup>184</sup>.  
1049 Similar to a GAN, GAN-like software contains a generator and a discriminator, but only one of  
1050 them is an artificial neural network.

Figure 1 | **The workflow for traditional and machine learning approaches in population genetic inference.** Traditional approaches build deterministic and stochastic models based on population genetics theory. Statistical inference can provide approaches to estimate parameters in stochastic models. Besides, population genetics theory can provide domain-specific knowledge when using machine learning algorithms. Currently, only supervised and unsupervised learning are applied in population genetic inference. Algorithms in supervised learning try to learn how to separate data and predict labels because training data are labelled (indicated by crosses and circles here), while algorithms in unsupervised learning try to learn probabilistic distributions to group data because no label is available in training data. Both supervised and unsupervised learning can use approaches based on deep learning or other algorithms. Simulation is supported by population genetics theory and used for validating different approaches. Also, it can provide simulated data for simulation-based approaches from statistical inference<sup>29</sup> and training data for machine learning algorithms<sup>15</sup>. Once an approach is validated with simulations, it can be applied to real data.

Figure 2 | **Common architectures and layers for artificial neural networks (ANNs).** a | ANNs contain  $m$  hidden layers with  $n$  neurons in each hidden layer. The arrows between nodes

represent weights or parameters learnt from data. ANNs can be feed-forward neural networks or recurrent neural networks (RNNs) depending on the information flow. An ANN usually contains three parts: an input layer, an output layer, and hidden layers. A typical convolutional neural network (CNN) combines  $q$  convolutional layers, activation layers, and pooling layers together. Their outputs then are passed into  $p$  fully connected layers. **b** | The input here is a biallelic genotype matrix. Different colours highlight variants from different positions. Neurons in the fully connected layer process all variants and assign different weights to the same allele, whereas neurons in the convolutional layer only process several variants depending on the kernel size and assign the same weight to the alleles from the same genome if the input is processed by the same kernel, leading to parameter sharing between different neurons. The kernel here is one-dimensional, and its weights are learnt from data. A max pooling layer outputs the maximum element from its inputs, enabling CNNs invariant to slight shifts in the input. **c** | RNNs can process biological sequences position by position while tracking information present in previous positions. **d** | Graph convolution is a generalized convolution because grid-like data are equivalent to graphs where a node (blue) connects with all its neighbouring nodes (red) inside a kernel. Then, convolution combines information with a node and those connecting with this node.

Figure 3 | **Deep generative models.** **a** | A restricted Boltzmann machine (RBM) contains only one visible (input) layer and only one hidden layer<sup>48</sup>. Information can flow between these two layers. Hence, RBMs are represented using undirected graphs, while other architectures are usually depicted using directed graphs with arrows. Genotype matrices can be regarded as images and passed into a deep generative model (DGM). Then the DGM can generate new data using the properties learnt from training data. In genotype matrices, different colours indicate different alleles at the same position. **b** | A variational autoencoder contains an encoder and a decoder<sup>48</sup>. The encoder converts the inputs into codes — latent representations. The decoder then tries to reconstruct the inputs similar to the original inputs and usually enforces the distribution of the code close to a normal distribution. The code can visualize population structure and the decoder can generate artificial genotypes<sup>83</sup>. **c** | A generative adversarial network (GAN) contains a generator and a discriminator<sup>48</sup>. The generator creates synthetic data and the discriminator tries to distinguish real and synthetic data. The generator and discriminator update their parameters in turn. The generator parameters are fixed and the discriminator parameters are updated for some iterations (blue arrows); then the discriminator parameters are fixed and the generator parameters are updated for some iterations (red arrows). **d** | GANs may have mode

collapse in which only subsets of the distribution of the training data are learnt, and mode dropping in which subsets of the distribution of the training data are forgotten.

Figure 4 | **Novel architectures and models.** **a** | A simplified transformer architecture is based on the vanilla transformer, which contains  $N$  encoders and decoders<sup>19</sup>. The encoder learns representation from training data through a series of layers, including self-attention layers, normalization layers, and feed-forward layers. Then the decoder uses this representation with training data to make predictions, for example, to predict whether a position in genomes is an introgressed variant (highlighted in red here) or not. Because transformers receive all inputs together, positional encoding is necessary if position information is important. **b** | Self-attention layers are the key components of transformers. The inputs are numeric vectors, such as genetic variants in genomes denoted with 0 and 1. Inside a self-attention layer, the inputs are converted into three components: query, key, and value. The outputs are generated using these three components through several operations, including matrix multiplication, transpose and softmax, which is a mathematical function that converts the elements of its inputs into values between 0 and 1 conditional on the sum of all the elements equal to 1. Here, the query weights, key weights, and value weights are parameters learnt from the data. Variants from different positions and their corresponding results are highlighted with different colours. **c** | Diffusion models contain forward and reverse processes. Here, the denoising diffusion probabilistic model<sup>202</sup> is illustrated as an example. Noise from a normal distribution is injected into the inputs, such as genotype matrices (different colours indicate different alleles at the same position), step by step during the forward diffusion process; then artificial neural networks can gradually recover the input by learning how to denoise during the reverse denoising process.

Figure 5 | **The implementation workflow.** To implement a deep learning model, researchers first define their problem and gather enough data before selecting an appropriate architecture. In population genetics, data could be either real data from biological sequences or simulated data from a given evolutionary model. The raw data can then undergo data pre-processing (Box 1). Typically, the processed dataset is split into training, validation, and test sets. The training set is further divided into smaller batches, which are randomly selected and fed into the model to iteratively update its parameters. When all batches have been utilized, an epoch is completed. Subsequently, the validation set is employed to monitor the model performance after each epoch. The choice of performance metrics relies on the specific problem and data type; a confusion matrix is depicted here as an example. If the model performance on the validation set is subpar,

researchers can explore different hyperparameter settings (grid search as an example here) and adjust the network architecture until satisfactory performance is achieved. Once the optimal model is obtained, it can be evaluated using the test set. If the performance remains unsatisfactory, researchers might consider collecting additional data, employing regularization techniques (weight decay as an example here), reducing the model complexity, or even replacing the goal with a less ambitious one. After experimenting these approaches, if the model still does not perform well, researchers may abandon deep learning; if the model demonstrates good performance on the test set, the workflow can conclude.

### **Box 1. Data preprocessing**

Typically, artificial neural networks (ANNs) require numeric inputs. Hence, non-numeric inputs, such as genome sequences, must be converted into numerical values using embedding or encoding techniques before being processed by ANNs. Encoding can transform non-numeric values into numbers (see the figure). A prevalent strategy for encoding genome sequences involves converting genome alignments with different nucleotides into genotype matrices containing 0 and 1 as elements, since population genetics usually assumes that a variant has only two allelic types. If missing values or multiple alleles must be considered, additional values such as  $-1$  can be incorporated into genotype matrices. Genotype matrices may be further filtered, such as by removing variants with many missing values, to ensure data quality. Another encoding approach utilizes summary statistics, such as the allele frequency spectrum<sup>45</sup>. In addition, one-hot encoding — a technique that converts each category variable into a distinct binary vector representation of several distinct categories, where 1 is placed at the position corresponding to the specific category and all other positions are filled with 0 — is common in machine learning, although it is not frequently applied in population genetics (but see ref. 61), which can handle multi-allelic data. However, encoding may not preserve the relationships between input features. Embedding can address this limitation and reduce input data dimensionality. Deep learning models can also be employed for dimensionality reduction to learn embeddings from data. A variational autoencoder was developed to learn embeddings in a hyperbolic space from genotype data. These embeddings can generate population trees for classifying individuals into different populations and can be used for studying genetic ancestry or simulating genotype data<sup>203</sup>. In addition to embeddings from genotype data, positional embedding could also be considered. For example, one study employed a separate fully connected layer to learn information from variant positions within genomes and merged position and genotype information for predictions<sup>63</sup>. Moreover, encoding and embedding can be used together<sup>61</sup>.

Furthermore, different neural network architectures may require specific techniques for pre-processing<sup>48</sup>. For example, if a convolutional neural network requires inputs with fixed size, techniques such as padding, which ensures that the inputs have the same size by adding additional elements around the original data, or width normalization can be considered<sup>75</sup>.

## Glossary

**Accuracy:** The ratio of the number of correct predictions to the total number of predictions in a dataset.

**Admixture:** The process in which genetic material from multiple populations merges into a single population.

**Adversarial attack:** A technique that slightly perturbs input data and causes machine learning models to make wrong predictions on such manipulated inputs with high confidence.

**Approximate Bayesian computation:** A statistical inference approach that employs simulation to estimate the posterior distribution of model parameters based on observed data when the exact posterior distribution is intractable, as motivated by Bayes Theorem.

**Backpropagation:** An optimization algorithm that recursively calculates gradients from the output layer to the input layer based on the chain rule from calculus for updating the parameters of an artificial neural network.

**Cross-entropy:** A metric that measures the performance of machine learning models for classification by comparing the similarity of two probability distributions.

**Decision tree:** A class of supervised learning algorithms that makes predictions by learning and organizing rules into a binary tree-like structure from data.

**Domain adaptation:** A technique that enables machine learning models trained with data from one domain (source domain) to be adapted and make accurate predictions on data from a different but related domain (target domain).

Dropout: A technique to regularize artificial neural networks by randomly deactivating neurons during training.

Early stopping: A technique to regularize artificial neural networks by stopping training before reaching the minimum loss on the training set.

F1 score: A metric that measures the performance of machine learning models for binary classification by calculating the harmonic mean of a given pair of precision and recall.

Game theory: A math discipline that studies strategies of interaction among rational players.

Gated recurrent unit: A unit similar to long short-term memory but with fewer parameters that improves the performance of recurrent neural networks on long sequences.

Grid search: A technique that optimizes hyperparameters by training and evaluating the model performance on combinations of a predefined set of hyperparameters.

Hidden Markov model: A class of generative models that processes sequential data by assuming the observed sequence is generated by a sequence of unobserved random variables independently transiting from the current state to the next state and not depending on all previous states.

Logistic regression: A type of regression that generates binary output ranging from 0 to 1 and can be viewed as a special type of artificial neural network composed by a neuron using a sigmoid activation function.

Long short-term memory: A unit that improves the performance of recurrent neural networks on long sequences by deciding whether information should be remembered or forgotten in the neural network along the sequence.

Loss function: A mathematical function that quantifies the difference (that is, loss) between the actual data and the predictions made by a machine learning model.

Markov chain: A sequence of random variables where the future state of a given step only depends on the current state and remains unaffected by all previous states.

Markov chain Monte Carlo: A statistical inference approach for estimating statistical model parameters by gradually approximating the probability distribution of model parameters by simulating a Markov chain of parameter values.

Maximum likelihood estimation: A statistical inference approach for estimating statistical model parameters by finding parameters that can maximize the probability of observed data.

Meta-learning: A class of machine learning algorithms that automates the learning process of machine learning algorithms for different tasks.

Mixture model: A probabilistic model that is generated from multiple atomic probability distributions.

Precision: The ratio of the number of correctly predicted instances to the total number of instances predicted as belonging to that class by the model in a dataset.

Principal component analysis: A technique that reduces the dimensionality of high-dimensional continuous data while keeping most of the information from the data by exploiting linear relationships among the features.

Random search: A technique that optimizes hyperparameters by training and evaluating the model performance on random combinations of hyperparameters from a predefined search space.

Recall: The ratio of the number of correctly predicted instances to the total number of instances actually belonging to that class in a dataset.

Rectified linear unit: A common activation function that returns the input value if the input value is larger than zero or returns zero otherwise.

Regression: A supervised learning task that makes quantitative predictions from input data.



Saliency map: A graphical representation that visualizes the contributions of each pixel in an image to the predictions made by an artificial neural network, revealing the regions of the image that significantly influence the decision-making process of the network.

Sequence-to-sequence learning: A machine learning task that involves training models to convert input sequences into corresponding output sequences, which is often employed in natural language processing for applications such as machine translation and speech recognition.

Simulated annealing approach: An optimization method that iteratively conducts probability solution updates proportional to the number of iterations already performed and the quality of the proposed solution to discover the global optimal solution.

Stochastic gradient descent: An optimization algorithm that iteratively updates parameters in machine learning models by randomly choosing data to calculate the gradients of the loss function.

Structured prediction: A supervised learning task that, unlike traditional classification or regression tasks which usually predict a single entity output, forecasts complex structures within the input data and generates outputs like sequences, trees, and graphs.

Summary statistic: A metric computed from genetic variants that is informative for an evolutionary parameter of interest.

Support: A set of values of a random variable for which the probabilities are greater than zero with a given probability distribution.

Tensor: A multidimensional array that is generalized from vectors and matrices for organizing high-dimensional data.

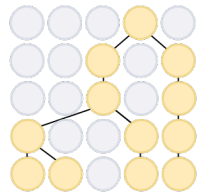
Transfer learning: A technique that allows reusing previously trained successful models for one task to other similar tasks.

Weight decay: A technique to regularize machine learning algorithms by penalizing large values in the model parameters through adding a term to the loss function that is proportional to the square of the magnitude of the parameters.

ACCEPTED MANUSCRIPT

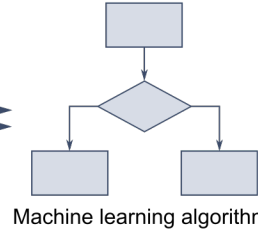
Theoretical support

Time before present ↑

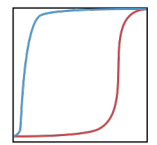


Population genetics theory

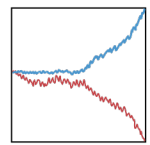
Domain-specific knowledge



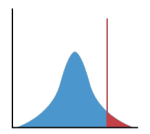
Machine learning algorithms



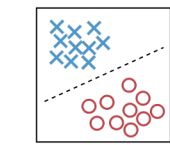
Deterministic model



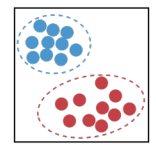
Stochastic model



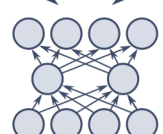
Statistical inference



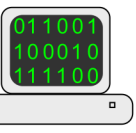
Supervised learning



Unsupervised learning



Deep learning

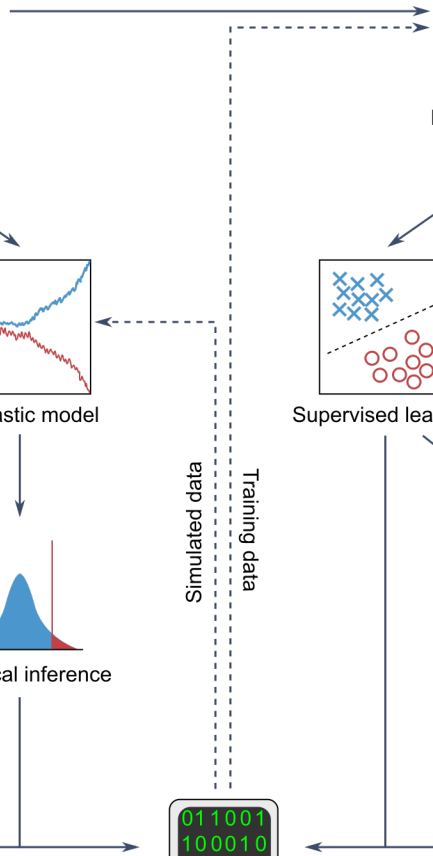
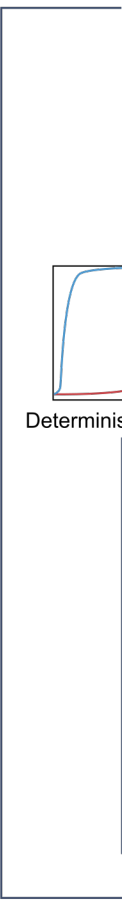


Simulation

Application



Real data



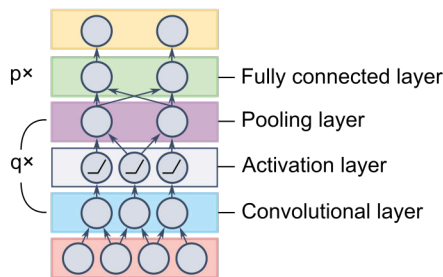
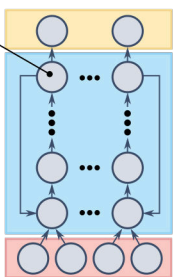
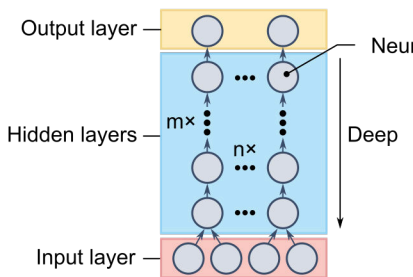
Validation

Validation

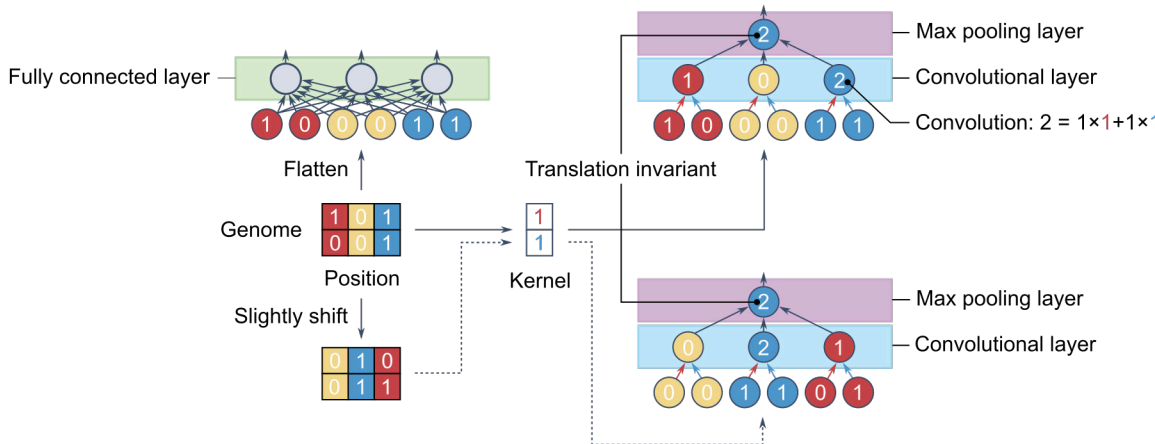
Feed-forward neural network

Recurrent neural network

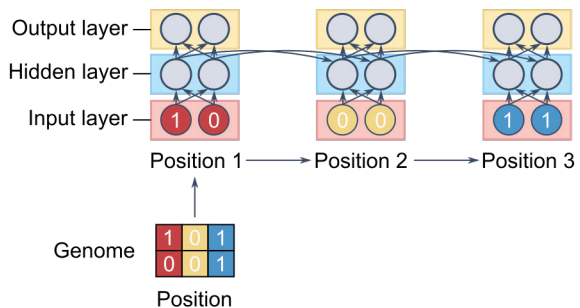
Convolutional neural network



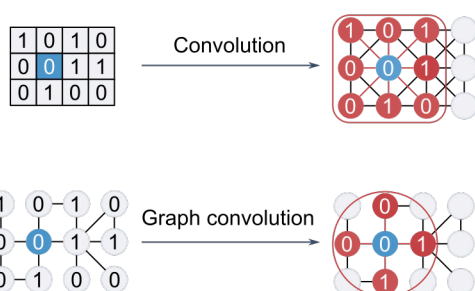
b

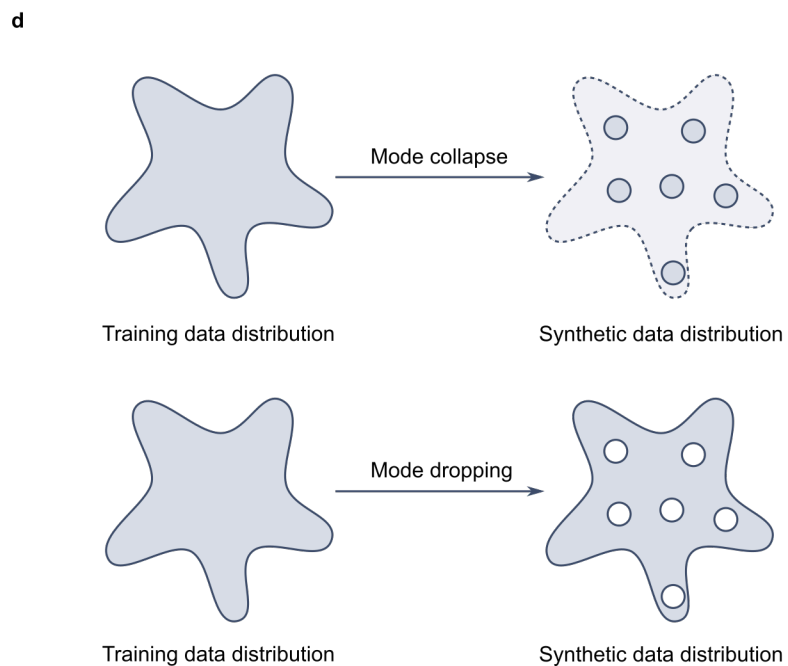
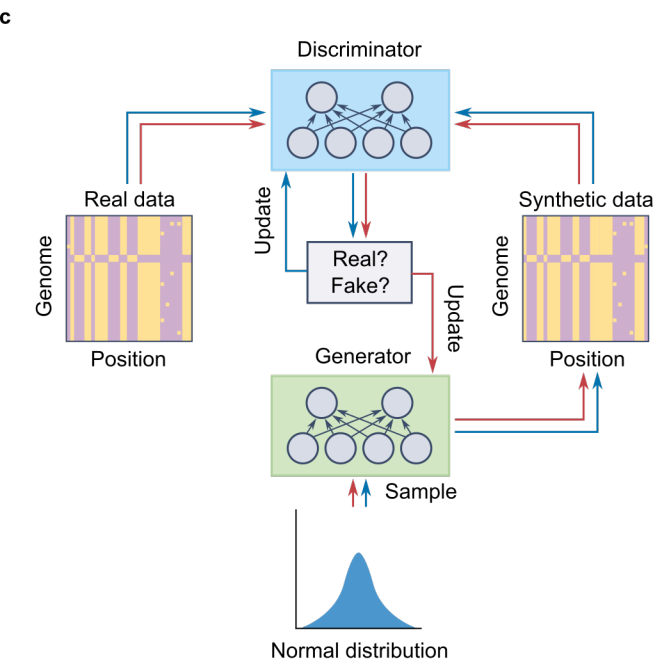
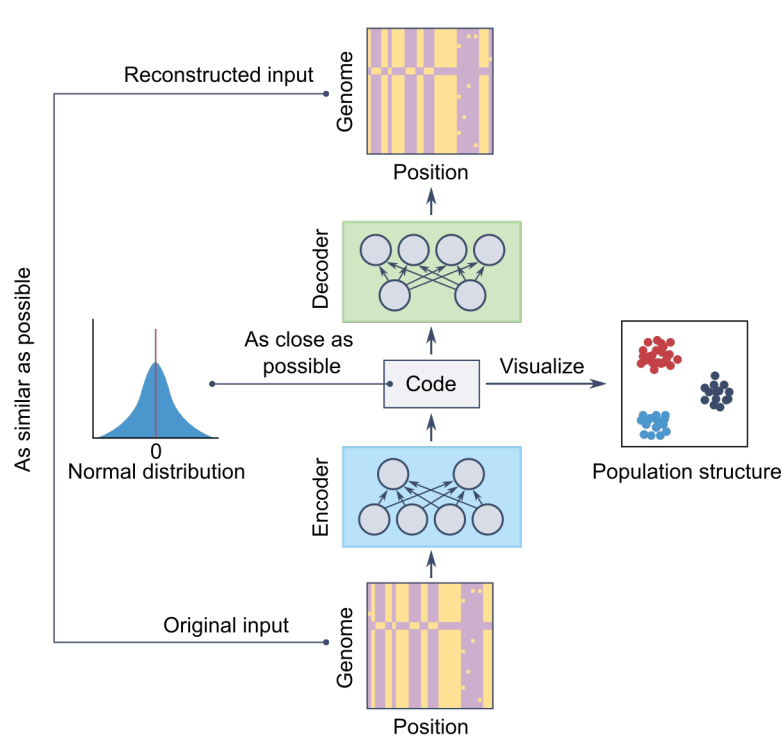
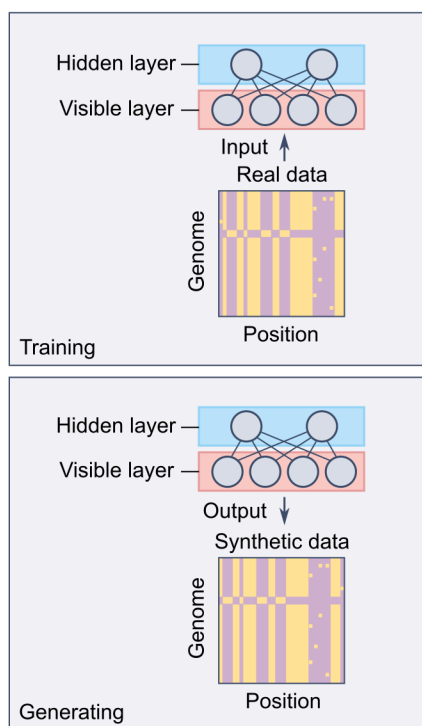


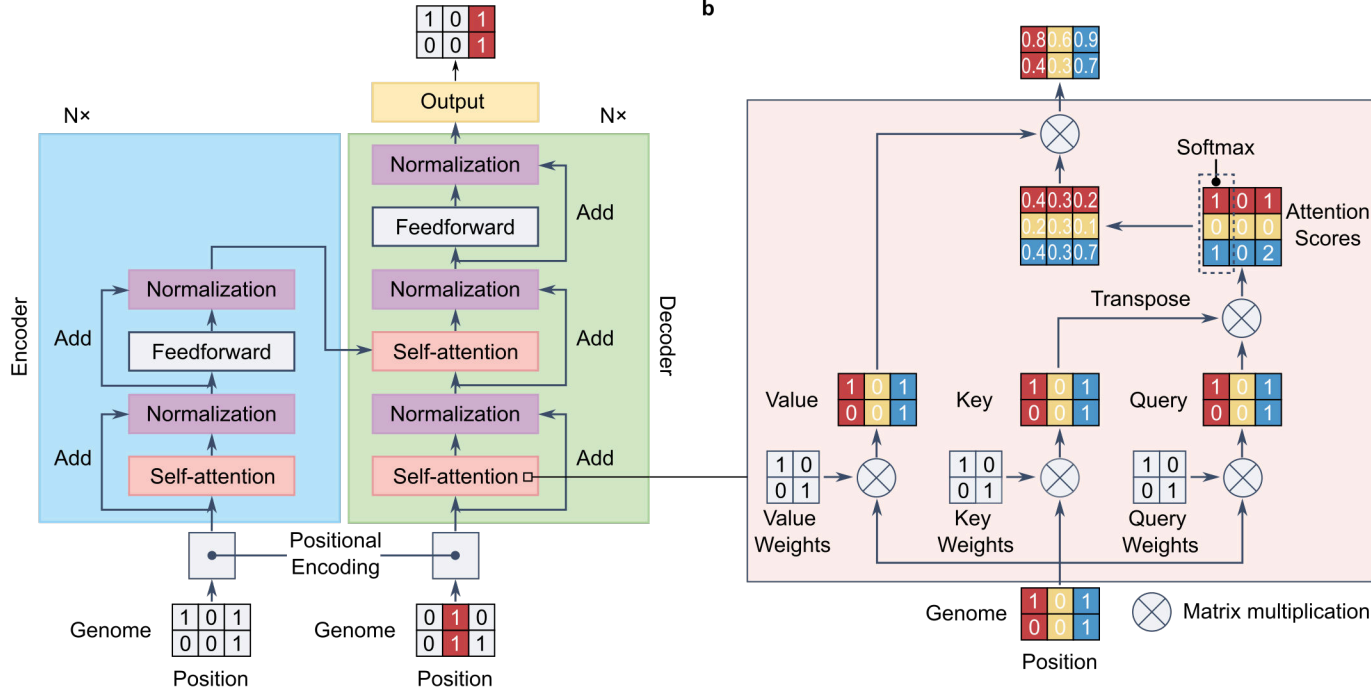
c



d







**c**

