# Deep learning for population size history inference: design, comparison and combination with approximate Bayesian computation

**Théophile Sanchez**[1*], **Jean Cury**[1], **Guillaume Charpiat**[1], **Flora Jay**[1*]

1. Laboratoire de Recherche en Informatique, CNRS UMR 8623, Université Paris-Sud, Université Paris-Saclay, Inria, Orsay, France

∗ Correspondence: `theophile.sanchez@inria.fr` and `flora.jay@lri.fr`

## Abstract

For the past decades, simulation-based likelihood-free inference methods have enabled to address numerous population genetics problems. As the richness and amount of simulated and real genetic data keep increasing, the field has a strong opportunity to tackle tasks that current methods hardly solve. However, high data dimensionality forces most methods to summarize large genomic datasets into a relatively small number of handcrafted features (summary statistics). Here we propose an alternative to summary statistics, based on the automatic extraction of relevant information using deep learning techniques. Specifically, we design artificial neural networks (ANNs) that take as input single nucleotide polymorphic sites (SNPs) found in individuals sampled from a single population and infer the past effective population size history. First, we provide guidelines to construct artificial neural networks that comply with the intrinsic properties of SNP data such as invariance to permutation of haplotypes, long scale interactions between SNPs and variable genomic length. Thanks to a Bayesian hyperparameter optimization procedure, we evaluate the performances of multiple networks and compare them to well established methods like Approximate Bayesian Computation (ABC). Even without the expert knowledge of summary statistics, our approach compares fairly well to an ABC based on handcrafted features. Furthermore we show that combining deep learning and ABC can improve performances while taking advantage of both frameworks. Finally, we apply our approach to reconstruct the effective population size history of cattle breed populations.

# 1   Introduction

In the past years, fields such as computer vision and natural language processing have shown impressive results thanks to the rise of deep learning methods. What makes these methods so powerful is not fully understood yet, but one key element is their ability to handle and exploit high dimensional structured data. Therefore, deep learning seems particularly suited to extract relevant information from genomic data, and has indeed been used for many tasks outside population genetics at first, such as prediction of protein binding sites, of phenotypes or of alternative splicing sites (Alipanahi et al., 2015, Jaganathan et al., 2019, Ma et al., 2018).

As genomic data becomes more and more available, it is possible to study the genetic variations within species or populations and investigate complex demographic histories including admixture events, population structure or size fluctuation through time, and this for many species. In fact, initiatives like the 1000 Genomes Project for human populations (Consortium et al., 2010) have been extended for better world coverage and data quality (Bergström et al., 2019, Consortium et al., 2015, Leitsalu et al., 2014, Mallick et al., 2016, Pagani et al., 2016) and opened up to many other species such as *Bos taurus* with the 1000 Bull Genomes Project (Daetwyler et al., 2014) or chimpanzees and gorillas with the Great Apes Genome Project (Prado-Martinez et al., 2013). Even for smaller scale studies, researchers have access to the whole genomes or high-density SNP data of numerous samples. These data collections can only be treated with inference methods able to scale to dozens or hundreds of individuals and large amount of genetic markers.

In this study, we propose several deep learning approaches for reconstructing the detailed histories of past effective population sizes from genetic polymorphism within a single population, a task considered difficult for various reasons. First, a present-day population, and even more so a sample of it, is one among many possible outcomes of a stochastic process depending on population sizes, mutation and recombination. Second, many other factors such as selective pressure, admixture events or population structure also shape the contemporary genetic diversity, which can blur the link between population size history and genetic data. As a result, the accuracy of the reconstruction and its level of resolution depend on the number of individuals available, the quality of the data and the methodology used. Nonetheless, in practice and under some simplifying assumptions, previous methods such as Bayesian skyline plots and their derivatives (Ho and Shapiro, 2011), sequential Markov coalescent (SMC) (PSMC, diCal and their derivatives (Li and Durbin, 2011, Sheehan et al., 2013)), Approximate Bayesian Computation (Boitard et al., 2016, Navascués et al., 2017) and SFS-based approaches (Bhaskar et al., 2015, Liu and Fu, 2015) have shown great results, supporting archaeological evidence and helping to understand species decline or expansion.

The study of genetic variation relies primarily on genotyping and sequencing data of very high dimensionality, which is a major difficulty for most inference methods. Some approaches, such as coalescent-HMMs methods (Spence et al., 2018), enable parameter inference using the full dataset by making simplifying assumptions on the underlying models. Despite impressive recent improvements, they still suffer some limitations: if a few of them can now process unphased data (Terhorst et al., 2017), scale to large sample size (Terhorst et al., 2017) or to complex models (Steinrücken et al., 2019), no method simultaneously address the three, and handling arbitrarily complex models remains untested (e.g.

2

52 models with more than three populations) or intractable (e.g. complex spatial models) (Spence et al., 2018). Hence, most

53 frameworks solving complex population genetic tasks do not rely on coalescent-HMMs and reduce the data dimension

54 with a pre-processing step during which the dataset is converted into a smaller set of statistics called summary statistics.

55 These statistics can then be used in likelihood and composite likelihood inference frameworks, when the model or

56 statistics are simple enough, or in simulation-based approaches. Among the latter, the widely used Approximate

57 Bayesian Computation (ABC) framework as well as several machine-learning algorithms, including SVM and random

58 forests, were able to tackle a variety of tasks such as demographic model selection and parameter inference (Excoffier

59 et al., 2013, Jay et al., 2019), detection of selection (Sugden et al., 2018, Tournebize et al., 2019) and introgression

60 (Schrider et al., 2018). The current trend when addressing complex tasks is to include a large number of summary

61 statistics inspired by population genetic theory in order to minimize the information loss. Summary statistics commonly

62 used are the site frequency spectrum (SFS) and its summaries (e.g. Tajima D), linkage disequilibrium (LD) and statistics

63 based on shared segments that are identical-by-state (IBS) or identical-by-descent (IBD) (Jay et al., 2019, Sheehan

64 and Song, 2016, Smith and Flaxman, 2019). However, they are not guaranteed to be sufficient and the inclusion of

65 numerous statistics can impact the performances of standard ABC, a problem known as curse of dimensionality (Blum,

66 2010). An active research topic in the ABC community is thus the development of methods addressing this curse of

67 dimensionality by (i) selecting the best subset of summary statistics according to some information-based criteria, or (ii)

68 integrating machine learning steps into ABC to handle a larger number of summary statistics (e.g. kernel methods,

69 random forests), or (iii) constructing summary statistics using linear and non-linear models based on candidate statistics

70 or on the original data when feasible (Aeschbacher et al., 2012, Blum et al., 2013, Fearnhead and Prangle, 2012, Jiang

71 et al., 2017, Nakagome et al., 2013, Raynal et al., 2018). In our study, we chose to use deep learning to bypass summary

72 statistics computation. The artificial neural network not only automatically computes its own summary statistics, also

73 called features in the deep learning community, directly from the genomic dataset, but also performs the inference task

74 to output effective population size predictions.

75 Deep learning, i.e. the use of artificial neural networks (ANN), has only very recently been used to tackle population

76 genetics questions. First, multilayer perceptron (MLP) were used to process large sets of summary statistics and to

77 predict jointly selective sweeps and simple demographic changes (Sheehan and Song, 2016), or to disentangle between

78 multiple scenarios of archaic introgression thanks to an additional ABC step (Lorente-Galdos et al., 2019, Mondal

79 et al., 2019). A second type of ANN, convolutional neural networks (CNN), were then applied to summary statistics

80 computed over 5Kb genomic regions in order to predict selective sweeps (Xue et al., 2019). A considerable shift

81 occurred when several studies applied ANN directly on genomic data instead of using summary statistic. Various

82 CNN architectures processing SNP matrices were proposed to infer recombination rates along the genome (Chan et al.,

83 2018, Flagel et al., 2018), selection (Flagel et al., 2018, Torada et al., 2019), introgression (Flagel et al., 2018) and

84 three-step population size histories (Flagel et al., 2018). The CNN implemented by Chan et al. (2018) and based on

85 Deep Sets (Zaheer et al., 2017) is invariant to haplotype permutation (i.e to the permutation of lines in the SNP matrix)

86 thanks to convolution filters that treat each haplotype in an identical way. The other approaches proposed instead to

sort haplotypes by similarity before processing them with filters sensitive to the haplotype order (Flagel et al., 2018, Torada et al., 2019). More recently, Recurrent Neural Networks (RNN) were applied to estimate the recombination rate (Adrion et al., 2019), and Generative Adversarial Networks (GAN) to learn the distribution of genomic datasets and to generate artificial genomes (Yelmen et al., 2019).

Convolution layers have been particularly efficient for large size data with spatial coherence such as images, exploiting the geometric structure of the image pixel grid. Instead of requiring as many weights as the input data size (like fully-connected layers in MLP), convolution layers take advantage of the spatial structure of the data, by defining spatially-small filters and applying them at each location along the input dimension (here, the SNP sequence). The result of each operation is ordered spatially according to the corresponding location in the input, to keep spatial coherence throughout the network. The filters of the first convolution layers have a small scope over the data input of the network, but adding layers on top of each other gives a larger scope to the last layers, allowing to handle long range interactions, in particular when combined with max-pooling layers. These interactions are important for demographic inference, e.g. they can allow the network to measure linkage disequilibrium and thus, the level of recombination.

Among the variety of developed ANN architectures, it is not straightforward to know which one is the most adapted to genomic data for a given population genetic task. In particular, this study aims at reconstructing detailed step-wise effective population size histories with 21 size parameters under an unknown recombination rate, a complex model with a fairly high dimensional parameter space compared to the population genetic task previously addressed with ANN. Hence we propose multiple networks, some of which are new and designed specifically for population genomics, and others that are more basic. We then apply a hyperparameter optimisation procedure (BOHB (Falkner et al., 2018)) to select the best architecture and hyperparameters. We investigate the performance of two MLPs, one using summary statistics and one using SNPs data of fixed length. We also compare two novel CNN based architectures, one with mixed convolution filter sizes over multiple individuals and another CNN that is adaptive to the genomic input size and invariant to the permutation of individuals or haplotypes. Both networks incorporate SNP data and their relative positions, a concept also developed in a different fashion by Flagel et al. (2018). In our last setup, we combine ABC and ANN by using the ANN predictions as summary statistics with the aim to benefit from both method advantages. Because no end-to-end deep learning approach for demographic inference had yet been compared to ABC or other traditional methods, we carefully benchmarked all these networks against variations of PopSizeABC, one of the highly performing methods for step-wise size inference that is based on ABC (Boitard et al., 2016). We also compare our architecture with CNNs developed for a related demographic task (Flagel et al., 2018). Finally we apply our approach to real genomes in order to reconstruct the size history of three cattle breeds.

## 2 Results

In this study, we introduce the first deep learning approaches for inferring detailed histories of effective population sizes using genomic data. Based on whole sequences of SNP data of multiple individuals from a single population, we aimed

120  to predict 21 population size parameters, each corresponding to a time step. Our method and the baseline frameworks

121  all relied on large-scale simulated datasets for which the true demographic parameters are known and drawn from prior

122  distributions of population sizes and recombination rates. For each drawn parameter set, we simulated 100 2Mb-long

123  genomic regions using msprime (Kelleher et al., 2016). Using this reference panel, we then trained methods based on

124  ABC, deep learning or a combination of both, to predict the demographic parameters (Figure 1). In this section, we will

125  give an overview of these methods as well as the hyperparameter optimisation procedure, followed by the evaluation of

126  their performance and the inference of the demographic history of three cattle breeds.

## 2.1  Baselines

128  We compare our approach to 5 baselines: an ABC approach and a MLP both using linkage disequilibrium and site

129  frequency spectrum as summary statistics, and another MLP, a "custom CNN" and a CNN from (Flagel et al., 2018), all

130  using genomic data directly. We evaluate four ABC methods (rejection, local linear regression, local ridge regression and

131  a single-hidden-layer neural network) with six tolerance rates (0.05, 0.1, 0.15, 0.2, 0.25 and 0.3). Their performances

132  represent a baseline for state-of-the-art methods. The second framework, a MLP processing the same summary statistics

133  as ABC, is similar in spirit to the one previously proposed for predicting 3 demographic parameters alongside selection

134  (Sheehan and Song, 2016). It is the baseline for methods using deep learning based on summary statistics. The third

135  and fourth baselines consist in a MLP and a CNN processing directly the first 400 SNPs of a 2Mb-long genomic region

136  instead of summary statistics. The MLP takes as input SNP data and locations flattened into a single vector. In this

137  configuration, the spatial information between SNPs, the link between the SNPs and their location, and the link between

138  alleles of the same individual are lost. The fourth baseline is a newly design CNN with rectangular shaped filters that

139  cover more than one haplotype (see Methods). This "custom CNN" is a first step towards an architecture tailored to raw

140  genomic data, because spatial information is preserved as for recent ANNs applied to population genetics (Chan et al.,

141  2018, Flagel et al., 2018, Torada et al., 2019), but also because asymmetrical filters of various sizes account for the

142  heterogeneous entities of axes (haplotype versus SNP, rather than pixel versus pixel). Finally we adapted and re-trained

143  four networks among the top-ranked CNNs proposed by Flagel et al. (2018) so that they could reconstruct a 21-epoch

144  model of instantaneous effective population size rather than the three-epoch model initially investigated by the authors,

145  and for practicability we called them *Flagel* CNNs.

## 2.2  Sequence Position Informed Deep Neural Architecture

147  We called our architecture SPIDNA, for Sequence Position Informed Deep Neural Architecture, and designed it to

148  comply to the principal features of SNP data: data heterogeneity (data includes genetic markers and their positions),

149  haplotype permutation invariance, long range dependencies between SNPs and variable number of SNPs.

### 2.2.1 Permutation invariance

One of the SNP matrix properties is its invariance to the permutation of haplotypes or genotypes. The same matrix with permuted rows contain the exact same information and should lead to the same predictions. Most summary statistics are already invariant to the haplotype order by definition. On the other hand, typical operations used in ANN such as rectangular filters and fully connected layers are not invariant, and consequently our two baseline ANNs do not respect this data feature. Here we implemented an architecture invariant by design, that stacks functions equivariant and invariant to row permutations (Lucas et al., 2018). It is a modification of the Deep Sets scheme (Zaheer et al., 2017) used for population genetics under the name "exchangeable networks" (Chan et al., 2018). In our study, the equivariant function is a convolutional layer with filters of size $1 \times a$, that treats each haplotype (row) independently and computes equivariant features, while the invariant function computes the mean of these features over the row dimension. The invariant function reduces the dimension of the data to one row, which is then concatenated to each equivariant row (Figure 2). Therefore the correlation between rows increases at each layer, which progressively transforms the equivariant input to an invariant output. However, the correlation increase should be moderate and progressive to avoid immediate loss of the information at the haplotype level. To promote this, we perform two independent normalizations, one over the output of the equivariant function and one over the input of the invariant function, and associate a correlation control parameter $\alpha$ that quantifies the contribution of the invariant function to the next layer, thus controlling the speed at which the correlation increases between rows.

### 2.2.2 Convolution networks to handle data with variable size

A major difficulty that arises with genomic data is that the number of SNP varies from one dataset to another, or from one genomic region to another, due to the stochasticity of biological and demographic processes (and of their corresponding genetic simulator). Therefore, we use convolution layers as they can handle data with variable size while keeping the number of network weights constant. A filter can be repeated more or fewer times to cover the whole input entering each layer, letting the network adapts itself to the data. Consequently, the output size of each convolution layer will vary depending on the input size. This prevents the use of fully connected layers directly after a convolution layer as it is often the case with CNNs. Instead, we use fully-connected layers only after operations independent of the input size and with a fixed output size, namely mean functions over the column and row dimensions (Figure 2).

Overall, we carefully designed an architecture accounting for invariance and adaptive specificities by stacking multiple equivariant blocks consisting in convolution layers with filters of size $1 \times 3$, invariant average pooling layers that are also adaptive to the number of SNPs and fully-connected layers that predict demographic parameters at each block (Figure 2). Contrary to our previous networks that include batch normalization of neuron activities, adaptive SPIDNA networks include instance normalization (see Methods).

6

## 2.3 Hyperparameter optimization

Compared to other machine learning methods, ANNs have a potentially infinite amount of hyperparameters when including for instance the number of layers, the number of neurons in each of them, the learning rate, weight decay or the batch size. Moreover, a run over a full dataset with enough epochs to reach convergence is time consuming for networks with a complex architecture defined by many learnable parameters. Therefore, the development of deep learning architectures often relies on the experience and intuition of the practitioner in a try-and-repeat process. Grid search and random search are two strategies for exploring the hyperparameter space uniformly. They are commonly used but are limited by the computing resources available. In our study, we used HpBandSter, a package that implements the HyperBand (Li et al., 2016) algorithm to run many hyperparameter trials on a smaller resource budget (i.e. few epochs) and runs the most promising trials on a greater budget. Combined with BOHB (Falkner et al., 2018), a Bayesian optimisation procedure that models the expected improvement of the joint hyperparameters, this method provides more guided and faster search of the hyperparameter space. At each step, BOHB draws a new combination of hyperparameter values to be tested according to the expected improvement and to a predefined prior. Here, we performed a search in a 5-dimensional space defined by: the type of architecture (architectures from our baselines and variations of SPIDNA architecture, based on 400 SNPs or the full number of SNPs), the learning rate, the weight decay, the batch size and the correlation control parameter $\alpha$. The search was performed for 3 budget steps and replicated 5 times, leading to a total of 83 successfully trained networks.

The configuration with the lowest loss generated by the hyperparameter optimization procedure used 400 SNPs with SPIDNA, batch normalization, a weight decay of $2.069 \cdot 10^{-2}$, a learning rate of $1.416 \cdot 10^{-2}$ and a batch size of 78 (Figure S1). Configurations with large batch sizes tended to have lower losses (Figure S1), which is expected as large batches provide a better approximation of the full training set gradient. However, a batch size too close to the training set size can lead to overfitting the training set. Here, we did not observe overfitting for any run when monitoring training and validation losses. The best configurations also tended to have low learning rates and weight decays (Figure S1). These low values slow down the convergence, but usually decrease the final prediction error if the budget is high enough for the network to reach convergence.

## 2.4 Comparison of the optimized architectures

For each architecture, we selected the best configuration obtained with the hyperparameter optimization procedure and trained it for a greater budget (i.e. 10 epochs), allowing an in-depth comparison. We found no strong decrease of prediction errors after this longer training compared to their counterparts with a $10^7$ budget ($10^7$ training SNP matrix, i.e. 5.57 epochs) (Figures 3 and S1).

We first compared the optimized neural networks to optimized ABC approaches based on predefined summary statistics. The prediction errors achieved by ABC using summary statistics ranged from 0.490 (ABC rejection, i.e. without correction) to 0.352 (ABC neural networks). The MLP network based on summary statistics performed comparatively

7

214   well (0.388). On the contrary the MLP based on raw data performed very poorly (0.690). All other networks based on

215   raw data outperformed this MLP, and most of them (all except SPIDNA instance normalization and SPIDNA instance

216   normalization adaptive) outperformed the ABC rejection (SPIDNA batch normalization, 0.453) or led to similar errors

217   (ranging from 0.485 to 0.489). The *Flagel* CNNs adapted from Flagel et al. (2018) that were not using dropout had

218   extreme overfitting issues (average training losses of 0.004 and 0.006 versus validation losses of 0.719 and 0.984). The

219   two other *Flagel* networks achieved prediction errors similar to SPIDNA (network based on the first 400 SNPs: 0.447;

220   network based on 1784 downsampled SNPs: 0.437), however they had 8 to 34 times more parameters than SPIDNA and

221   despite the dropout procedure they still suffered from overfitting (validation over training loss ratio of 2.50 and 1.39 for

222   *Flagel* CNNs versus 0.89, 1.01 and 0.97 for SPIDNA which indeed sometimes performed better on the validation set).

223   Lastly, we evaluated two methods that combine deep learning and ABC, by considering the features automatically

224   computed by a network as summary statistics for ABC (Jiang et al., 2017). When using only the predictions of SPIDNA

225   as input to ABC with correction (linear regression, ridge regression or neural network), we improved greatly SPIDNA's

226   performances and obtained errors similar to the ABC based on predefined summary statistics (0.363 compared to 0.352).

227   When using both SPIDNA predictions and predefined summary statistics as input to the ABC algorithm we decreased

228   further the prediction errors (0.335).

229   We further illustrated the performances of SPIDNA on a subset of demographic scenarios that were previously

230   investigated (Boitard et al., 2016) (Figure 4). The method correctly reconstructed histories of constant size, expansion

231   and decline, as SPIDNA predictions from 100 independent genomic regions (black boxplots) approximately followed

232   the real population size trend and magnitude. The true parameters were always included in the 90% credible intervals

233   (light blue envelopes) predicted by SPIDNA combined with ABC without predefined summary statistics and, for

234   most cases, in the 50% credible intervals (dark blue). Both methods also correctly reconstructed a complex history

235   encompassing an expansion interrupted by a bottleneck and followed by a constant size (see Figure 4 'Bottleneck').

236   However, they were unable to correctly estimate the parameters of a very complex 'Zigzag' history except for its initial

237   growth period and instead reconstructed a smoother history with values intermediate to the lower and higher population

238   sizes (see Figure 4 'Zigzag'). This confirmed the smoothing behavior identified previously for ABC and MSMC

239   (Boitard et al., 2016). Finally, similarly to ABC on predefined summary statistics (Boitard et al., 2016), SPIDNA

240   predictions of very recent population sizes were slightly biased toward the center of the prior distribution, however

241   combining SPIDNA with ABC tended to correct this bias in most cases.

## 2.5   SPIDNA infers the decline of effective population size of cattle

243   We inferred the effective population size history of three breeds of cattle (Angus, Fleckvieh and Holstein) based on the

244   same 75 individuals studied by Boitard et al. (2016) and sampled by the 1000 Bull Genomes Project (Figure 5). The best

245   ABC and SPIDNA configurations both infer a large ancestral effective population size and a decline for the past 70,000

246   years. However, SPIDNA reports higher recent population size (Angus:11,334, Holstein:12,311, Fleckvieh:13,579) than

247   ABC (Angus:361, Holstein:552, Fleckvieh:1,329). Interestingly, SPIDNA infers the same population sizes for all three

8

breeds before 10 thousand years ago. This is in agreement with the estimation of the beginning of the domestication (Zeder, 2008). SPIDNA combined with ABC also reconstructed a decline after domestication but estimated larger population sizes for the last 30,000 years than SPIDNA alone. In addition Angus had the largest recent population and Fleckvieh the smallest in contrary to the two previous methods. Finally SPIDNA combined with ABC identified an episode of smooth decline and recovery of the population size preceding the domestication in the putative ancestral species (between 400,000 and 30,000 years ago).

# 3 Methods

## 3.1 Simulated data and summary statistics

All methods compared in this study are trained in a supervised fashion, and thus require genetic data of numerous populations under various demographic scenarios. We defined the demographic parameters by following similar rules as Boitard et al. (2016): $I = 21$ time windows $[t_i, t_{i+1}]$ were defined from present to ancient periods with $t_i = (exp(log(1 + aT)i/(I - 1)) - 1)/a$ generations, $i$ going from 0 to $I - 1$, $T = 130,000$, $a = 0.06$ and $t_I = +\infty$. By increasing exponentially the time windows as we go further in the past, we obtain more detailed scenarios for recent times. The time windows are identical for all scenarios to focus on the population size parameters. Each demographic scenario is generated by drawing a first population size $N_0$ between 10 and 100,000 from a uniform distribution and corresponds to the most recent time window $t_0$. The population sizes of the next time windows follow $N_i = N_{i-1} \times 10^{\beta}$ with $\beta$ sampled uniformly between -1 and 1. $\beta$ is redrawn if it gives a population size out of $]10, 100000[$. 50000 scenarios are drawn from this prior distribution and are passed to the msprime coalescent simulator version 0.6.1 (Kelleher et al., 2016) to simulate 100 independent 2Mb-long segments of 50 haploid individuals. Each simulation gives a SNP matrix $X$ of size $M$ haplotypes $\times S$ SNP sites and a location vector of size $S$ that features the relative locations between SNPs in the corresponding DNA sequence. Ancestral and derived alleles are encoded with 0 and 1. The mutation rate is set to $10^{-8}$ and the recombination rate is sampled uniformly between $10^{-9}$ and $10^{-8}$ for each scenario. The real cattle dataset has 4357 SNPs on average, thus we remove the scenarios with lower than 400 SNPs as we consider them outside of our prior. That reduced the dataset to 18461 scenarios, with 100 examples each, that we split into a validation set of 500 scenarios and a training set with the remaining 17961 scenarios.

For each group of 100 segments corresponding to one scenario, we computed the site frequency spectrum and the linkage disequilibrium as a function of the distance between SNPs averaged over 19 distance bins for a total of 87 summary statistics. Our python script is partly based on the scikit-allel python module (Miles et al., 2019).

## 3.2 Baseline specifications

**Approximate Bayesian Computation** We compared ABC with the simple rejection procedure (i.e. no correction) and three correction methods implemented in the R package 'abc' (Csilléry et al., 2012): local linear regression, ridge regression and non-linear regression based on a single-hidden-layer neural network. Settings were set to default except

9

280 for the tolerance rate. ABC was run on (a) Predefined summary statistics, (b) SPIDNA outputs (i.e. automatically

281 computed summary statistics), or (c) a combination of predefined summary statistics and SPIDNA outputs.

282 **Multi-Layer-Perceptron Networks** Our MLP based on summary statistics has 3 hidden layers, ReLU activation

283 functions, and uses batch normalisation and the adam optimizer. As in Sheehan and Song (2016), the hidden layers

284 have respectively 25, 25, and 10 neurons. It takes 35 summary statistics as input. This network and all the following

285 ones output 21 demographic parameters and are trained with a regular L2 loss function unless stated otherwise. This

286 MLP has a total of 2,986 trainable parameters. Our second MLP is based on 'raw' genomic data and takes a matrix of

287 50 haplotypes (rows) for 400 SNPs (columns) and its associated position vector as input. Its hidden layers respectively

288 have 20, 20, and 10 neurons, which gives it 408,981 trainable parameters.

289 **Custom CNN** Our convolutional neural network takes as input the same matrix of 400 SNPs and has 2-dimension

290 filters of various shapes. The first layer consists of 5 kernels with rectangular shape ($2 \times 2$, $5 \times 4$, $3 \times 8$, $2 \times 10$, $20 \times 1$)

291 applied to the SNP matrix $X$. Each kernel creates 50 filters, which amounts to 250 feature maps after the first layer.

292 The position vector $d$ is treated by the 5 associated kernel shapes ($1 \times 2$, $1 \times 4$, $1 \times 8$, $1 \times 10$, $1 \times 1$) with 20 filters

293 each, making 100 filters in total. The results of the first convolutional layer are then concatenated so that the second

294 convolutional layer will couple information from $X$ and $d$ in a way that emphasizes the original location of the SNPs

295 along the genome. The outputs of this second layer are then combined and go through 5 convolutional layers and 2

296 fully connected layers. Adding convolutional layers one after each other allows our network to combine patterns and

297 reduce the size of the data without adding too many weights to our model. This network has a total of 131,731 trainable

298 parameters.

299 *Flagel* **network** We reused the code associated with the repository of the first paper using a CNN for demographic

300 inference (Flagel et al., 2018) and adapted it to our dataset and task. We trained the network with the exact same

301 architecture as the one published, except that we changed the last layer to allow the prediction of our 21 population

302 size parameters. We parametrized the network with the set of hyperparameters leading to the best performance in the

303 previous work for two different types of SNP encoding (0/255 or -1/1). It is noteworthy that the actual encoding in their

304 code is 0/-1 and not 0/255, thus we kept the same encoding to be able to compare the performances. The networks were

305 trained with the same procedure of 10 epochs with early stopping in case of no progression of the loss after 3 epochs.

306 The batch size is 200. The input data had 50 haplotypes and either 400 SNPs as processed by our custom CNN or we

307 downsampled the data to one every ten SNPs as done in the original work, leading to 1784 wide input SNP matrices.

308 This size corresponds to the tenth of the biggest SNP matrix in our dataset. Smaller simulations are padded with zeros.

309 All parameters can be found in table S1.

310 ### 3.3 SPIDNA architecture

311 For our permutation invariant architecture, we designed three variations. The first one uses batch normalization, after

312 each convolution layer, and therefore takes as input a fixed number of 400 SNPs, similarly to two of the baselines.

The second one is invariant to the number of SNPs and uses instance normalization, after each convolution layer, to normalize layer inputs per-data instead of per-batch (for the batch normalization). The last variation is also invariant to the number of SNPs, but uses two separate instance normalization steps, as well as the correlation control parameter $\alpha$.

Except for these differences, these three variations have the same architecture, represented in Figure 2. At each step $i$ of the network, we consider that the data has four dimensions $B_i \times M_i \times S_i \times F_i$, $B$ being the batch dimension, $M$ the row dimension (also the haplotype/genotype dimension before the first layer), $S$ the column dimension (also the SNP dimension before the first layer) and $F$ the feature dimension (only one feature before the first layer). A first convolution layer of 50 $1 \times 3$ filters is applied to the SNP matrix, and another convolution layer of 50 $1 \times 3$ filters is applied to the position vector and repeated $M$ times. The results of the two convolutions have now the same dimensions and are concatenated along the feature dimension. The resulting tensor is then passed to seven blocks put end to end, each one involving an equivariant function and an invariant function. The equivariant function $\psi$ is a convolutional layer of 50 $1 \times 3$ filters that outputs a tensor of size $B_{i-1} \times M_{i-1} \times (S_{i-1} - 2) \times F_{i-1}/2$. The result of the equivariant function is then passed to the invariant function $\rho$, which is the mean over the dimension $M$. Thus $\rho(\phi(X_{i-1}))$ has size $B_{i-1} \times (S_{i-1} - 2) \times F_{i-1}/2$, which is repeated $M$ times to maintain the same dimension as $\phi(X_{i-1})$. Then $\rho(\phi(X_{i-1}))$ and $\phi(X_{i-1})$ are concatenated over the feature dimension. Finally, max-pooling filters of dimension $1 \times 2$ are applied, and the result is passed to the next block. In parallel, each block computes the average over the column dimension $S$ of the 21 first features of $\rho(\phi(X_{i-1}))$ that are then passed to a fully-connected layer with 21 outputs. The predictions of each block are summed.

### 3.4   From batch normalization to instance normalization

Network weight initialization is a difficult task that can lead to vanishing or exploding gradient when not carefully done and associated with a poor learning rate (Bengio et al., 1994, Glorot and Bengio, 2010). Most initialization schemes try to force the outputs of each layer to follow some distribution assuming normalized input data. Batch normalization solves this problem by normalizing layer outputs over the whole batch during training and computing a running mean and variance for the evaluation steps. We used this type of normalization for our networks that take as input a fixed number of SNPs. For the networks invariant to the number of SNPs, we could not concatenate all batch data into the same tensor because of their varying sizes. Therefore, we use instance normalization, which computes both mean and variance over the feature dimension.

### 3.5   Hyperparameter optimization

**ABC**   The tolerance rates ranged from 0.05 to 0.3 by step of 0.5 and were optimized for 12 ABC algorithms independently (4 correction methods $\times$ 3 types of inputs: predefined summary statistics, SPIDNA outputs or both).

**MLP using summary statistics**   As the training time for this algorithm was short, we optimized its hyperparameters with a random search by drawing 27 configurations from uniform distributions and trained a network for each

345 configuration during 6 epochs. The batch size was drawn between 10 and 100, learning rate between $5 \cdot 10^{-5}$ and

346 $1 \cdot 10^{-2}$ and weight decay between $5 \cdot 10^{-5}$ and $1 \cdot 10^{-2}$.

347 **ANNs over SNP matrices**    Batch size, learning rate, weight decay and the architecture choice were optimized using

348 the Bayesian optimization procedure BOHB (Falkner et al., 2018) with the same hyperparameter prior distributions

349 than the MLP using summary statistics, and with a uniform distribution over the architecture choices. For SPIDNA

350 architectures that controlled correlation, we added the control parameter $\alpha$ to the Bayesian optimization procedure with

351 a log-uniform prior between 0.5 and 1.

### 3.6   Cattle breed data

353 We inferred the demographic history of Angus, Fleckvieh and Holstein cattle breeds using the data set of 25 sequenced

354 individuals from the 1,000 genome bull project (Daetwyler et al., 2014) that was analysed by (Boitard et al., 2016). As

355 the data of real cattle sequence is prone to phasing and sequencing errors, we converted the real data from haplotype to

356 genotype with a minimum allele frequency (maf) of 0.2, as suggested by Boitard et al. (2016) and applied the same

357 treatment to the simulated training set. We split the real data of each breed into 1213 segments of 2Mb and removed

358 segments comprising centromeres. Then we trained ABC and SPIDNA with the best hyperparameter configurations on

359 the modified simulated data and performed the inference.

## 4   Discussion

361 In this paper, we introduced a deep learning approach to infer the detailed size history of a single population directly

362 from genomic data given an unknown recombination rate. This consisted in inferring jointly 21 population size

363 parameters. We not only increased the complexity of the demographic model with respect to previous works such as

364 Flagel et al. (2018), but also compared the performances of our architecture to other methods including ABC, and

365 applied our approach to real data sets. We found that our approach compared competitively with one of the best to

366 date approaches, with the added advantage of not relying on summary statistics. Finally, we reconstructed the effective

367 population size fluctuations of three cattle breeds and confirmed that they all had similar sizes when they were part

368 of the same ancestral species *Bos taurus* and underwent a decline likely linked to their domestication, although the

369 estimated strength of this decline depended on the inference method.

### 4.1   On the practicability and importance of architecture design

371 When applying deep learning techniques, the design of the neural network architecture is critical, as poor design can lead

372 to a lack of expressive power, information loss, underfitting, overfitting, or unnecessary complications that slow down

373 the training process. The recent history of successes in Computer Vision consists in architecture improvements, leading

374 to performance gaps (e.g. from MLP to LeNet, AlexNet, VGG, Inception and ResNet (He et al., 2016, Krizhevsky

375 et al., 2012, LeCun et al., 1998, Simonyan and Zisserman, 2014, Szegedy et al., 2017)). But these successes have been

12

built incrementally by relatively small changes over the last years, involving a large number of studies, researchers and challenges. Indeed finding the best architecture suited for a task is hard and time-consuming given the wide range of possibilities. Therefore, automating architecture and hyperparameter choice is an important challenge that can yield benefit to smaller fields such as population genetics. In our study, the Bayesian hyperparmeter optimisation procedure allowed to test multiple networks thanks to a better usage of the computational power available by giving more budget to the most promising ANN architectures and hyperparameters. This procedure could be extended to hyperparameters that further describe the architecture of the network such as the number and type of layers, of neurons, the type of non-linearity or the topology. Thanks to this procedure we investigated a series of architectures, starting from the simple multi-layer fully-connected network (MLP) and moving on to more complex architectures, and exhibited the link between design and performances.

To interpret the results and compare them, let us first note that in Figure 3, a 0 error means perfect prediction, while an error of 1 means that no information is extracted from the input. Indeed, a function outputting always the same value, for all samples, can at best predict the average target value over the dataset, in which case the error is the standard deviation over the dataset of the value to predict, which is normalized to 1 in our setup.

Processing the SNP and position matrices with a MLP led to high prediction errors, especially for recent population sizes. This is not surprising, since genomic information is encoded as a simple list of values, where the order has no meaning from the MLP point of view, which then cannot exploit information given by the data structure. In summary, an MLP configuration has several drawbacks: (i) the number of network parameters to estimate is high; (ii) the MLP can only retrieve the geometry of the data through training, with no guarantee that it will learn the spatial structure of the genome (i.e. the column order and distance between SNPs) or distinguish from which individual comes each SNP. In spite of all these hindrances, the MLP still performed far better than random guesses or constant prediction (31% better).

On the contrary, CNN layers process input elements by groups, allowing close SNPs to be processed together. This feature, combined with the stacking of layers in CNN, helps the network to construct features dependent of the SNPs proximity. Important summary statistics used in ABC or other inference methods such as linkage disequilibrium can potentially be easily expressed by such CNN. Illustrating the difficulty of undertaking the resolution of a new task with previously designed networks, the adapted *Flagel* CNNs (Flagel et al., 2018) suffered overfitting. Hence we proposed several novel convolutional architectures, tailored to genetic data. We first developed a CNN with 2D filters that could have different shapes, i.e. mixed kernel sizes but also non symmetrical masks. There is indeed no rational behind considering square masks only as is usually done in computer vision to describe pixel neighborhoods, as rows and columns in our case correspond to different entities (individual or phased haplotype versus markers). Using varied mask shapes (e.g., $2 \times 2$, $5 \times 4$ or $3 \times 8$) helps our custom CNN to learn features of various patterns, potentially mimicking different types of summary statistics ("vertical" masks integrate over individuals, enabling the computation of allele frequencies at a SNP, while "horizontal" ones integrate over SNPs, as IBS or IBD sharing tract length does). Such mixed size filters have proved useful in the Computer Vision literature also, under the name of Inception architectures

13

(Szegedy et al., 2017); they allow to extract and mix different kinds of information from multiple scales within the same layer. The large gap in performance between a simple MLP and this custom CNN confirms the importance of such considerations. A natural extension would be to integrate this feature into SPIDNA, our permutation-invariant architecture.

## 4.2 Novel architectures tailored to genomic data

### 4.2.1 Invariance to haplotype permutation

The order in which simulated haplotypes are arranged in a SNP matrix has no meaning. Although the custom CNN network above cannot be guaranteed to be exactly invariant to the haplotype order, it can approximately learn this specificity. To avoid wasting training time to learn that there is no information in the row order, it has been proposed to systematically sort the haplotypes according to a predefined rule (Flagel et al., 2018, Torada et al., 2019). Because there is no ordering in high dimensional space that is stable with respect to perturbations (Qi et al., 2017), we chose yet another alternative and enforced our network to be permutation-invariant by design. Permutation-invariant networks, or exchangeable networks, were successfully applied in population genetics by Chan et al. (2018) for inferring local recombination, but our architecture is different in that the invariant operations are performed at each block, enabling both individual equivariant features and global invariant features to contribute to the next layer. It has been proven that this type of architecture provides universal approximation of permutation-invariant functions (?). Among our permutation-invariant architectures, the best one (SPIDNA using batch normalization) had a smaller prediction error than our custom CNN. However, it is not clear whether this improvement is directly linked to its built-in permutation-invariance property, or to other differences between the two networks. Controlling the speed to invariance thanks to the parameter $\alpha$ improved the performance of the instance normalization SPIDNA, but not significantly the performance of the instance normalization adaptive SPIDNA (see table S2).

Importantly, SPIDNA adapts to the number of individuals, which is an advantageous property compared to many methods relying on summary statistics. SPIDNA can be trained on data sets having similar or varying sample sizes, and, once trained, it can be directly applied to a dataset of sample size unobserved during training.

### 4.2.2 Automatic adaptation to the number of SNPs

The two networks designed to be adaptive to the number of SNPs have the advantage to be applicable to genetic data of any length, to the opposite of networks specific to a particular number of SNPs, which transform the data with padding or compression, or are retrained for different lengths, or take as input portions of larger sequences. Our two SPIDNA adaptive networks show results close to the best of non-adaptive versions, though slightly worse (0.487 versus 0.453, see table S2), although the difference disappears when SPIDNA is combined with ABC (0.363 versus 0.363). This small performance drop is likely due to differences in normalization rather than to the adaptive feature. Indeed, the best non-adaptive SPIDNA uses batch normalization while the adaptive versions use instance normalization, because of implementation limitations, as there is currently no implementation of batch normalization for batches with inputs of

14

444 mixed sizes. Nonetheless, SPIDNA networks with instance normalization had much better performances when using all
445 SNPs rather than the 400 first SNPs only, which suggests that adaptability is a useful feature (see table S2).

446 Our adaptive architecture provides an alternative to data compression based on computer vision algorithms: since
447 compression is not optimized for the task of interest, it could induce information loss by reducing data prematurely.
448 Note indeed that the success of deep learning in computer vision lies precisely in the replacement of ad-hoc data
449 encoding (e.g. bag of words or SIFT (Lowe, 2004)) by ones that can be optimized. It is also an alternative to padding, a
450 technique that consists in completing the SNP and position data at the edges so that they all match the biggest simulated
451 SNP matrix; it is left to the neural network to guess where the real genetic data stops and where padding starts. As such
452 it may make the task more difficult, given that the SNP matrix size is highly variable between different demographic
453 histories and some examples would contain more padding values that actual genetic information. RNN are also a natural
454 alternative, though they induce an unequal contribution of SNPs to the final result, depending on their ordering along
455 the sequence. They were very recently proven to be useful to predict local recombination rate along the genome (Adrion
456 et al., 2019) and future works should investigate whether this scales up to a global characteristic and to a different task.

### 4.3 Advantages and challenges of deep learning

458 Alongside the quality of deep learning to automatically extract informative features from high dimensional data, artificial
459 neural networks are also very flexible. For instance, they can be used for transfer learning, that is, a network trained for
460 a specific task can be reused for another one by only modifying the last layers (e.g. a network trained for population
461 size history inference could be reused for classification between demographic scenarios) (Pan and Yang, 2009). The
462 new network will benefit from the embedding already learned for the previous task, improving error and learning time.
463 We also highlight that, as for most ABC methods, the parameters are inferred jointly, a major point as the common
464 population genetics model parameters almost never have an independent impact on shaping genetic diversity. We noted
465 that for highly fluctuating population sizes, SPIDNA estimated smooth histories. Smoothing can be seen as a good
466 byproduct and was for example achieved by SMC++ using a spline regularization scheme (Terhorst et al., 2017).

#### 4.3.1 Combining deep learning and Approximate Bayesian Computation to approximate the posterior distribution

469 We found that adding an ABC step to the deep learning approach increased its performance. This ABC step takes as
470 input the demographic parameters predicted by SPIDNA instead of the usual summary statistics. This strategy was
471 proposed by Jiang et al. (2017) that showed that a deep neural network could approximate the parameter posterior
472 means, which are desirable summary statistics for ABC. It was applied under the name of ABC-DL in two population
473 genetics studies for performing model selection, however both papers relied on the joint SFS as predefined candidate
474 summary statistics (Lorente-Galdos et al., 2019, Mondal et al., 2019). Here, we are taking advantage of both the deep
475 architecture to bypass summary statistics and the Bayesian framework to refine the prediction and approximate the
476 posterior distribution. Additionally, we showed that applying ABC to SPIDNA predictions combined with precomputed

15

477  summary statistics led to an error 3.6% smaller than the one of a regular ABC and 6.7% smaller than SPIDNA. This

478  indicates that the information retrieved by SPIDNA does not completely overlap the one encoded into the predefined

479  summary statistics but is not completely orthogonal either. It is a first step towards understanding and interpreting the

480  artificial neural networks currently used in population genetics, a major challenge that the deep learning field currently

481  faces for many of its applications (Gilpin et al., 2018) and that has not yet been investigated in our community.

### 4.3.2   Application to real data

483  Applying a method trained on simulated data to a real dataset can be a difficult task. Here we show that the estimated

484  effective population sizes of the three cattle breeds were qualitatively similar across the different methods used. All of

485  them were able to recover the large ancestral population size shared by the three breeds, followed by its decline after

486  domestication. However, the methods produced size estimates that were quantitatively different, notably in the strength

487  of the decline and the recent population sizes. For quality reasons, inference was done using genotypes pruned of low

488  frequency alleles rather than haplotypes. The architecture and hyperparameters were optimized based on simulated

489  haplotypes, and the network was trained on simulated genotypes. It is possible that an architecture designed with

490  a new hyperparameter optimization procedure calibrated for filtered genotypes would decrease SPIDNA error rate.

491  However, the discrepancy between ABC and SPIDNA reconstructions in the last 10,000 years might also be due to the

492  sensitivity of ANNs to overfitting and to mispecifications in the model generating training data. For example, decrease

493  in performances due to demographic mispecification has already been shown for selection inference based on ANNs

494  (Torada et al., 2019). In our case, model mispecification arises because cattle breeds are subjected to strong artificial

495  selection pressures, with few males contributing to the next generations, which is a clear violation of the coalescent

496  assumptions underlying our training simulated set. In addition, errors or missing information in real data were not

497  modelled in the training set, a procedure that can improve ABC performances when using multiple summary statistics

498  such as haplotype length statistics (Jay et al., 2019). When comparing performances on training and validation sets, we

499  found that our architectures were not overfitting, contrary to the architectures directly adapted from Flagel et al. (2018).

500  Yet it is possible that the features automatically constructed by ANNs are more sensitive to a gap between real and

501  simulated data (e.g. unmodelled errors and artificial selection) than an ABC method based on SFS and LD statistics.

502  Systematically testing and improving the robustness of ANNs trained on simulations is a great challenge for the coming

503  years.

### Conclusion

505  We addressed a challenging task of population genetics, that is, reconstructing effective population size through time.

506  We showed that this demographic inference could be done for unknown recombination rates. We are confident that

507  a network exchangeable and adaptive to the input size is a promising architecture for future lines of works for other

508  population genetics tasks, as it could prevent premature loss of information and favor learning new features rather

509  than known haplotype invariance. As for now, co-estimating multiple processes remains a hard task, and inference

16

510 is mostly done under a simplifying assumption, e.g. inferring selection or recombination under a fixed demographic

511 scenario, or inferring step-wise population size for a single population (but see MSMC and MSMC-IM for an extension

512 to two populations with interactions (Schiffels and Durbin, 2014, Wang et al., 2019)). The success of ABC and

513 simulation-based methods are partly due to their convenience to include complex models via simulations. Here we

514 showed, for the first time, that a well designed artificial neural network is capable of retrieving information about

515 fluctuating effective population size and that it can also be combined with existing summary statistics. Additionally,

516 recent studies demonstrated very good performances of other networks for solving tasks such as detecting introgression

517 and selection. For the above reasons, and because extracting information automatically should lead to the identification

518 of features that disentangle processes hardly distinguishable, we are hopeful that future robust networks trained on

519 complex simulations could help solving jointly some of these tasks. Finally, we provided (i) a tool for users willing

520 to infer population size history of any species that can be applied to phased or unphased genomes (available from

521 `https://gitlab.inria.fr/ml_genetics/public/DNA_DNA`); (ii) new exchangeable network architectures, some

522 of which have the promising feature of being adaptive to input size ; (iii) guidelines for future developers on building

523 architectures and hyper-optimization to facilitate the development of new artificial neural networks for population

524 genomics.

## Acknowledgments

## Data Availability Statement

532 Data sharing is not applicable to this article as no new data were created or analyzed in this study.

## References

534 Jeffrey R Adrion, Jared G Galloway, and Andrew D Kern. Inferring the landscape of recombination using recurrent

535     neural networks. *bioRxiv*, page 662247, 2019.

536 Simon Aeschbacher, Mark A Beaumont, and Andreas Futschik. A novel approach for choosing summary statistics in

537     approximate bayesian computation. *Genetics*, 192(3):1027–1047, 2012.

538 Babak Alipanahi, Andrew Delong, Matthew T Weirauch, and Brendan J Frey. Predicting the sequence specificities of

539     dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 33(8):831, 2015.

17

Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *Trans. Neur. Netw.*, 5(2):157–166, March 1994. ISSN 1045-9227. doi: 10.1109/72.279181. URL https://doi.org/10.1109/72.279181.

Anders Bergström, Shane A McCarthy, Ruoyun Hui, Mohamed A Almarri, Qasim Ayub, Petr Danecek, Yuan Chen, Sabine Felkel, Pille Hallast, Jack Kamm, et al. Insights into human genetic variation and population history from 929 diverse genomes. *bioRxiv*, page 674986, 2019.

Anand Bhaskar, YX Rachel Wang, and Yun S Song. Efficient inference of population size histories and locus-specific mutation rates from large-sample genomic variation data. *Genome research*, 25(2):268–279, 2015.

Michael GB Blum. Approximate bayesian computation: a nonparametric perspective. *Journal of the American Statistical Association*, 105(491):1178–1187, 2010.

Michael GB Blum, Maria Antonieta Nunes, Dennis Prangle, Scott A Sisson, et al. A comparative review of dimension reduction methods in approximate bayesian computation. *Statistical Science*, 28(2):189–208, 2013.

Simon Boitard, Willy Rodriguez, Flora Jay, Stefano Mona, and Frédéric Austerlitz. Inferring population size history from large samples of genome-wide molecular data-an approximate bayesian computation approach. *PLoS genetics*, 12(3):e1005877, 2016.

Jeffrey Chan, Valerio Perrone, Jeffrey Spence, Paul Jenkins, Sara Mathieson, and Yun Song. A likelihood-free inference framework for population genetic data using exchangeable neural networks. In *Advances in Neural Information Processing Systems*, pages 8594–8605, 2018.

1000 Genomes Project Consortium et al. A map of human genome variation from population-scale sequencing. *Nature*, 467(7319):1061, 2010.

1000 Genomes Project Consortium et al. A global reference for human genetic variation. *Nature*, 526(7571):68, 2015.

Katalin Csilléry, Olivier François, and Michael GB Blum. abc: an r package for approximate bayesian computation (abc). *Methods in ecology and evolution*, 3(3):475–479, 2012.

Hans D Daetwyler, Aurélien Capitan, Hubert Pausch, Paul Stothard, Rianne Van Binsbergen, Rasmus F Brøndum, Xiaoping Liao, Anis Djari, Sabrina C Rodriguez, Cécile Grohs, et al. Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nature genetics*, 46(8):858, 2014.

Laurent Excoffier, Isabelle Dupanloup, Emilia Huerta-Sánchez, Vitor C Sousa, and Matthieu Foll. Robust demographic inference from genomic and snp data. *PLoS genetics*, 9(10):e1003905, 2013.

Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and efficient hyperparameter optimization at scale. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1437–1446, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/falkner18a.html.

18

572 Paul Fearnhead and Dennis Prangle. Constructing summary statistics for approximate bayesian computation: semi-
573     automatic approximate bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Method-*
574     *ology)*, 74(3):419–474, 2012.

575 Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The unreasonable effectiveness of convolutional neural networks
576     in population genetic inference. *Molecular biology and evolution*, 36(2):220–238, 2018.

577 Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining explanations:
578     An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and*
579     *advanced analytics (DSAA)*, pages 80–89. IEEE, 2018.

580 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In
581     *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS'10). Society for*
582     *Artificial Intelligence and Statistics*, 2010.

583 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings*
584     *of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

585 Simon YW Ho and Beth Shapiro. Skyline-plot methods for estimating demographic history from nucleotide sequences.
586     *Molecular ecology resources*, 11(3):423–434, 2011.

587 Kishore Jaganathan, Sofia Kyriazopoulou Panagiotopoulou, Jeremy F McRae, Siavash Fazel Darbandi, David Knowles,
588     Yang I Li, Jack A Kosmicki, Juan Arbelaez, Wenwu Cui, Grace B Schwartz, et al. Predicting splicing from primary
589     sequence with deep learning. *Cell*, 176(3):535–548, 2019.

590 Flora Jay, Simon Boitard, and Frédéric Austerlitz. An abc method for whole-genome sequence data: inferring paleolithic
591     and neolithic human expansions. *Molecular biology and evolution*, 36(7):1565–1579, 2019.

592 Bai Jiang, Tung-yu Wu, Charles Zheng, and Wing H Wong. Learning summary statistic for approximate bayesian
593     computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.

594 Jerome Kelleher, Alison M Etheridge, and Gilean McVean. Efficient coalescent simulation and genealogical analysis
595     for large sample sizes. *PLoS Comput Biol*, 12(5):1–22, 05 2016. doi: 10.1371/journal.pcbi.1004842. URL
596     `http://dx.doi.org/10.1371%2Fjournal.pcbi.1004842`.

597 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural
598     networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

599 Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document
600     recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

601 Liis Leitsalu, Toomas Haller, Tõnu Esko, Mari-Liis Tammesoo, Helene Alavere, Harold Snieder, Markus Perola,
602     Pauline C Ng, Reedik Mägi, Lili Milani, et al. Cohort profile: Estonian biobank of the estonian genome center,
603     university of tartu. *International journal of epidemiology*, 44(4):1137–1147, 2014.

Heng Li and Richard Durbin. Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357):493, 2011.

Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *arXiv preprint arXiv:1603.06560*, 2016.

Xiaoming Liu and Yun-Xin Fu. Exploring population size changes using snp frequency spectra. *Nature genetics*, 47(5): 555, 2015.

Belen Lorente-Galdos, Oscar Lao, Gerard Serra-Vidal, Gabriel Santpere, Lukas FK Kuderna, Lara R Arauna, Karima Fadhlaoui-Zid, Ville N Pimenoff, Himla Soodyall, Pierre Zalloua, et al. Whole-genome sequence analysis of a pan african set of samples reveals archaic gene flow from an extinct basal population of modern humans into sub-saharan populations. *Genome biology*, 20(1):77, 2019.

David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60 (2):91–110, 2004.

Thomas Lucas, Corentin Tallec, Yann Ollivier, and Jakob Verbeek. Mixed batches and symmetric discriminators for GAN training. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2844–2853, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/lucas18a.html.

Wenlong Ma, Zhixu Qiu, Jie Song, Jiajia Li, Qian Cheng, Jingjing Zhai, and Chuang Ma. A deep convolutional neural network approach for predicting phenotypes from genotypes. *Planta*, 248(5):1307–1318, 2018.

Swapan Mallick, Heng Li, Mark Lipson, Iain Mathieson, Melissa Gymrek, Fernando Racimo, Mengyao Zhao, Niru Chennagiri, Susanne Nordenfelt, Arti Tandon, et al. The simons genome diversity project: 300 genomes from 142 diverse populations. *Nature*, 538(7624):201, 2016.

Alistair Miles, Peter Ralph, Summer Rae, and Rahul Pisupati. cggh/scikit-allel: v1.2.1, June 2019. URL https://doi.org/10.5281/zenodo.3238280.

Mayukh Mondal, Jaume Bertranpetit, and Oscar Lao. Approximate bayesian computation with deep learning supports a third archaic introgression in asia and oceania. *Nature communications*, 10(1):246, 2019.

Shigeki Nakagome, Kenji Fukumizu, and Shuhei Mano. Kernel approximate bayesian computation in population genetic inferences. *Statistical applications in genetics and molecular biology*, 12(6):667–678, 2013.

Miguel Navascués, Raphaël Leblois, and Concetta Burgarella. Demographic inference through approximate-bayesian-computation skyline plots. *PeerJ*, 5:e3530, 2017.

Luca Pagani, Daniel John Lawson, Evelyn Jagoda, Alexander Mörseburg, Anders Eriksson, Mario Mitt, Florian Clemente, Georgi Hudjashov, Michael DeGiorgio, Lauri Saag, et al. Genomic analyses inform on migration events during the peopling of eurasia. *Nature*, 538(7624):238, 2016.

636  Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*,
637      22(10):1345–1359, 2009.

638  Javier Prado-Martinez, Peter H Sudmant, Jeffrey M Kidd, Heng Li, Joanna L Kelley, Belen Lorente-Galdos, Krishna R
639      Veeramah, August E Woerner, Timothy D O'Connor, Gabriel Santpere, et al. Great ape genetic diversity and
640      population history. *Nature*, 499(7459):471, 2013.

641  Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification
642      and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
643      652–660, 2017.

644  Louis Raynal, Jean-Michel Marin, Pierre Pudlo, Mathieu Ribatet, Christian P Robert, and Arnaud Estoup. Abc random
645      forests for bayesian parameter inference. *Bioinformatics*, 35(10):1720–1728, 2018.

646  Stephan Schiffels and Richard Durbin. Inferring human population size and separation history from multiple genome
647      sequences. *Nature genetics*, 46(8):919, 2014.

648  Daniel R Schrider, Julien Ayroles, Daniel R Matute, and Andrew D Kern. Supervised machine learning reveals
649      introgressed loci in the genomes of drosophila simulans and d. sechellia. *PLoS genetics*, 14(4):e1007341, 2018.

650  Sara Sheehan and Yun S Song. Deep learning for population genetic inference. *PLoS computational biology*, 12(3):
651      e1004845, 2016.

652  Sara Sheehan, Kelley Harris, and Yun S Song. Estimating variable effective population sizes from multiple genomes: a
653      sequentially markov conditional sampling distribution approach. *Genetics*, 194(3):647–662, 2013.

654  Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv
655      preprint arXiv:1409.1556*, 2014.

656  Chris CR Smith and Samuel M Flaxman. Leveraging whole genome sequencing data for demographic inference with
657      approximate bayesian computation. *Molecular ecology resources*, 2019.

658  Jeffrey P Spence, Matthias Steinrücken, Jonathan Terhorst, and Yun S Song. Inference of population history using
659      coalescent hmms: review and outlook. *Current opinion in genetics & development*, 53:70–76, 2018.

660  Matthias Steinrücken, Jack Kamm, Jeffrey P Spence, and Yun S Song. Inference of complex population histories using
661      whole-genome sequences from multiple populations. *Proceedings of the National Academy of Sciences*, 116(34):
662      17115–17120, 2019.

663  Lauren Alpert Sugden, Elizabeth G Atkinson, Annie P Fischer, Stephen Rong, Brenna M Henn, and Sohini Ramachan-
664      dran. Localization of adaptive variants in human genomes using averaged one-dependence estimation. *Nature
665      communications*, 9(1):703, 2018.

666  Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the
667      impact of residual connections on learning. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

668  Jonathan Terhorst, John A Kamm, and Yun S Song. Robust and scalable inference of population history from hundreds
669    of unphased whole genomes. *Nature genetics*, 49(2):303, 2017.

670  Luis Torada, Lucrezia Lorenzon, Alice Beddis, Ulas Isildak, Linda Pattini, Sara Mathieson, and Matteo Fumagalli.
671    Imagene: a convolutional neural network to quantify natural selection from genomic data. *BMC bioinformatics*, 20
672    (9):337, 2019.

673  Rémi Tournebize, Valérie Poncet, Mattias Jakobsson, Yves Vigouroux, and Stéphanie Manel. Mcswan: A joint site
674    frequency spectrum method to detect and date selective sweeps across multiple population genomes. *Molecular*
675    *ecology resources*, 19(1):283–295, 2019.

676  Ke Wang, Iain Mathieson, Jared O'Connell, and Stephan Schiffels. Tracking human population structure through time
677    from whole genome sequences. *bioRxiv*, page 585265, 2019.

678  Alexander T Xue, Daniel R Schrider, Andrew D Kern, Ag1000G Consortium, et al. Discovery of ongoing selective
679    sweeps within anopheles mosquito populations using deep learning. *bioRxiv*, page 589069, 2019.

680  Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec, Francesco Montinaro, Cyril
681    Furtlehner, Luca Pagani, and Flora Jay. Creating artificial human genomes using generative models. *bioRxiv*, page
682    769091, 2019.

683  Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan R Salakhutdinov, and Alexander J
684    Smola. Deep sets. In *Advances in neural information processing systems*, pages 3391–3401, 2017.

685  Melinda A Zeder. Domestication and early agriculture in the mediterranean basin: Origins, diffusion, and impact.
686    *Proceedings of the national Academy of Sciences*, 105(33):11597–11604, 2008.
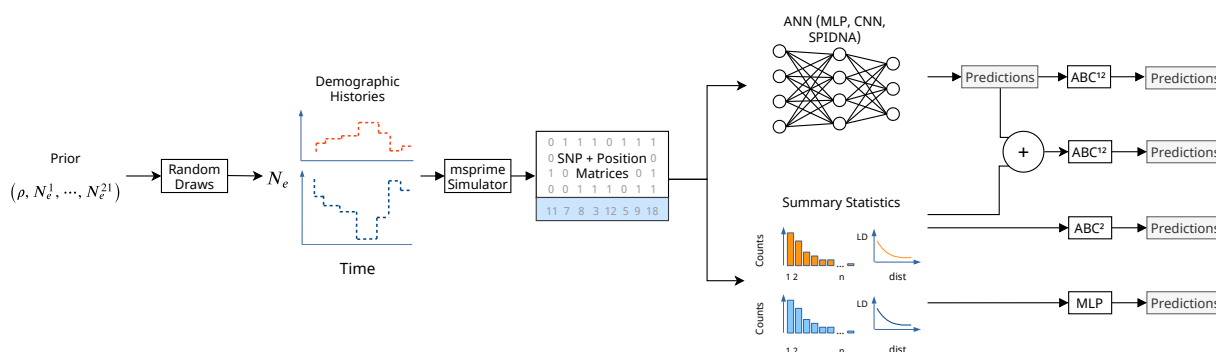
687  # Figures and Tables



Figure 1: Overview of the methods compared in this study. Demographic histories are drawn from a prior distribution on 21 population sizes $N_e^i$ and one recombination rate $\rho$, and are used to generate SNP matrices with msprime. Two types of summary statistics are computed from these simulations (SFS and LD). The predictions (outputs) made by different kind of ANNs (MLP, custom CNN and SPIDNA architecture) are compared to an MLP using the summary statistics and to ABC using either the summary statistics, SPIDNA outputs or both. [1] ANN outputs used are the predictions made by the version of SPIDNA with the lowest prediction error. [2] ABC without correction, with linear regression, ridge regression or a single layer neural network are compared.
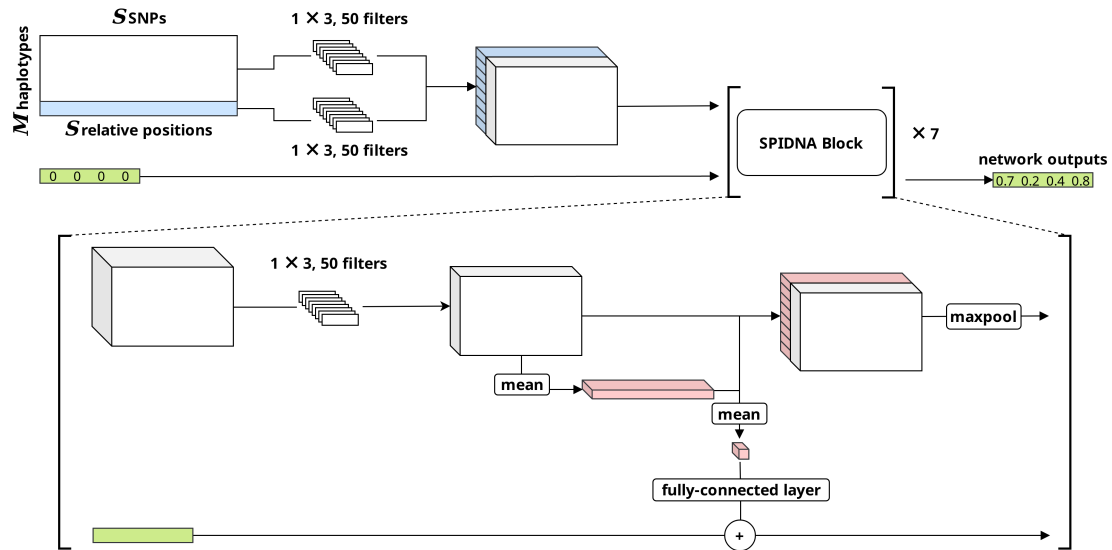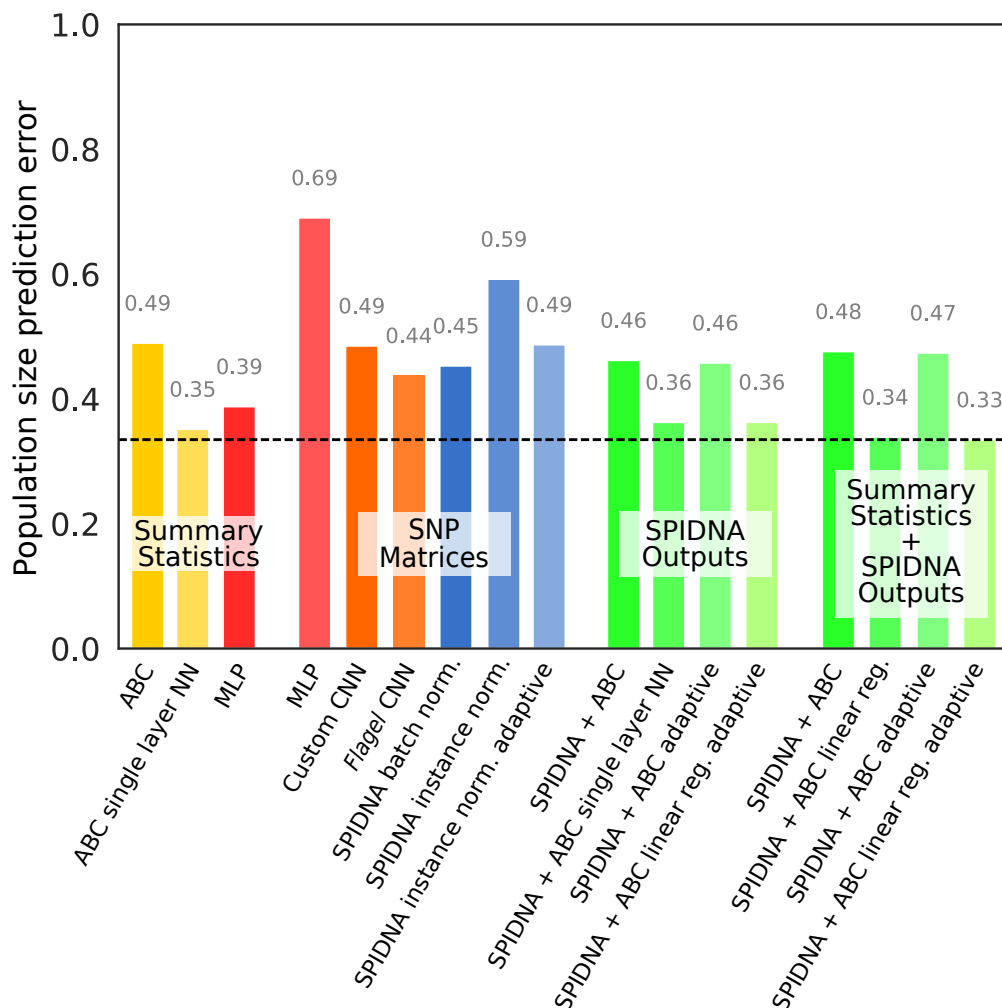
22

Figure 2: Schematic of SPIDNA architecture.

Figure 3: Prediction errors of the best run of each method after the hyperparameter optimization. The best configurations of each ANN (MLP, custom CNN and SPIDNA) have been retrained for 10 epochs. Traditional ABC methods are depicted in yellow, deep MLPs and CNNs in red and orange, SPIDNA ANNs in blue, combinations of ANNs and ABC in green. Methods are grouped into 4 families: "Summary statistics" (processed by ABC or ANN), "SNP matrices" (processed by ANN), "SPIDNA outputs" (processed by ABC, no summary statistic used), "Summary statistics and SPIDNA outputs" (processed by ABC).
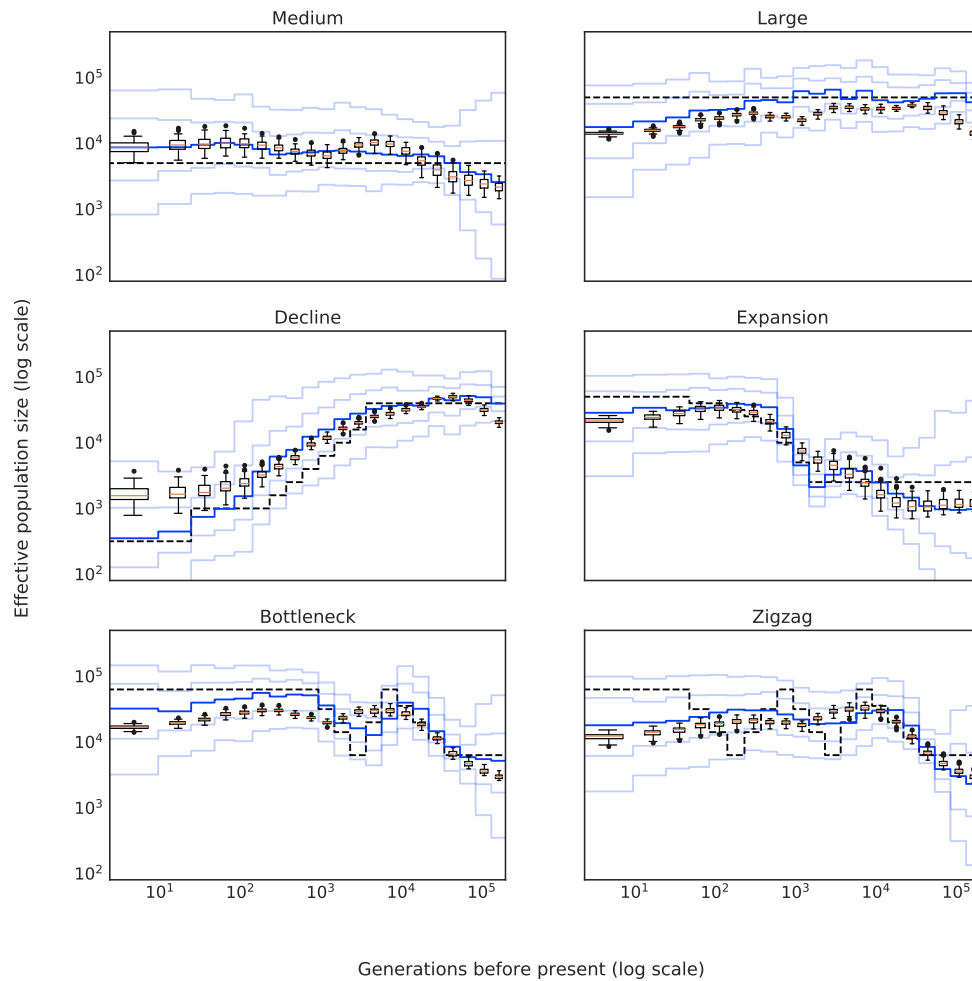
Figure 4: Predictions of SPIDNA and ABC using SPIDNA outputs, for six predefined scenarios (dashed black lines). 100 replicates were simulated for each scenario. Boxplots show the dispersion of SPIDNA predictions (over replicates). For each history inferred by SPIDNA combined with ABC, we display the posterior median (plain blue line), the 50% credible interval (dark blue) and the 90% credible interval (light blue).
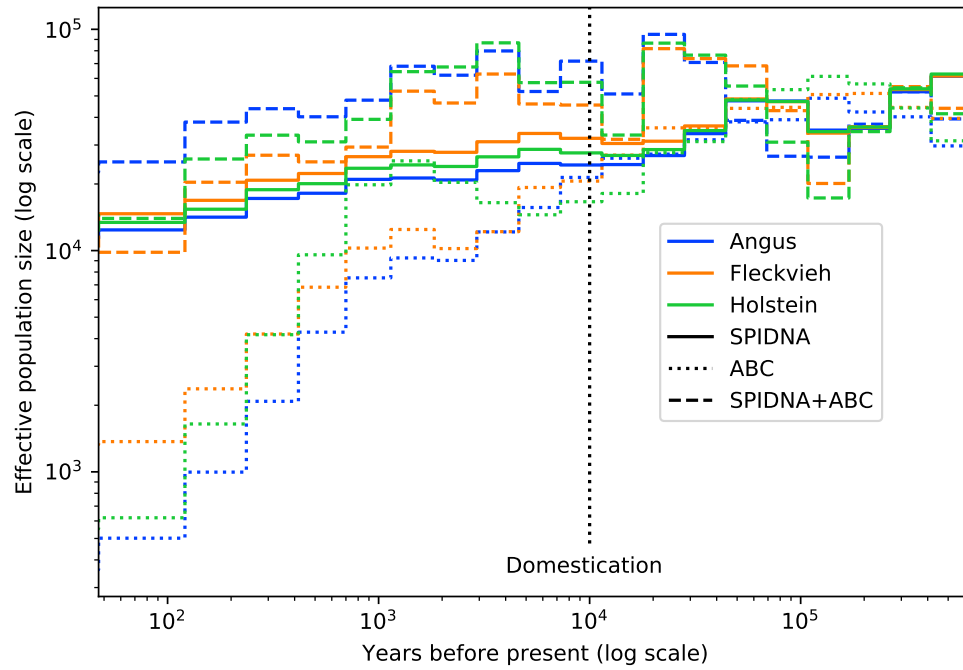
Figure 5: Effective population size of three cattle breeds inferred by ABC (dotted lines), by the best SPIDNA architecture, SPIDNA batch normalization (dashed lines), and by ABC based on SPIDNA outputs (dashed lines). Domestication is estimated to have occurred 10 000 years ago (vertical dotted line).
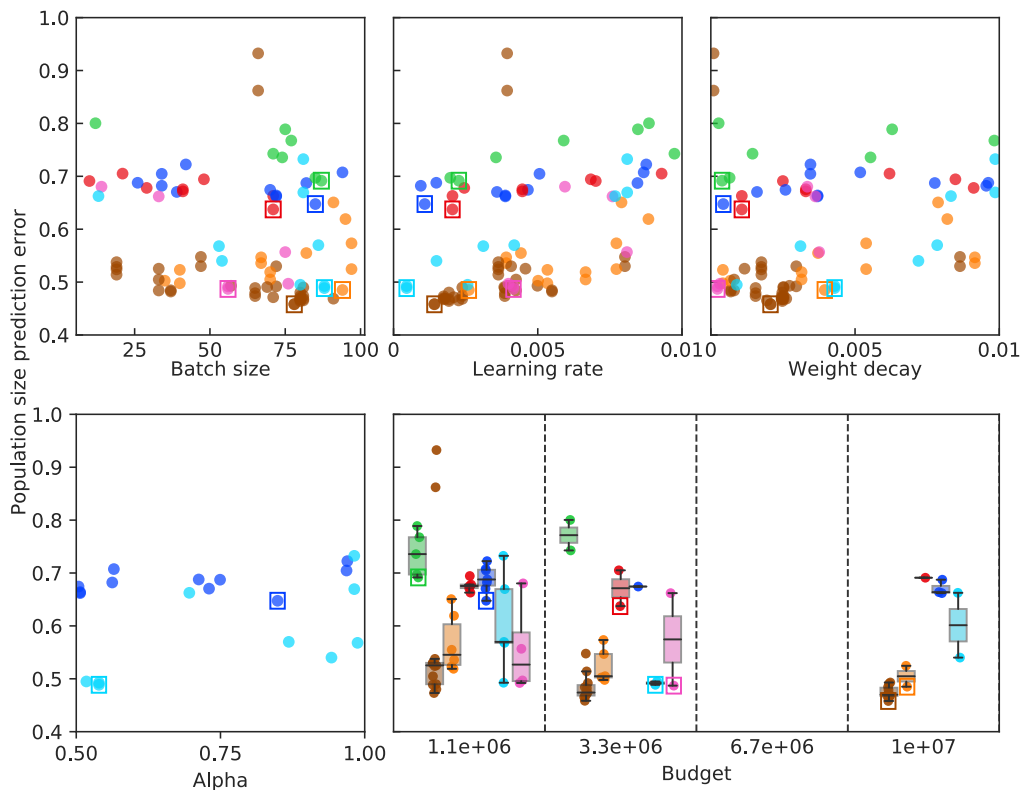
26

688 ## Supplementary materials



Figure S1: Population size prediction error for each run of the hyperparameter optimization procedure. X-axes indicate the hyperparameter (batch size, learning rate, weight decay and alpha) or budget values, and colors indicate the type of network used for the run (MLP, custom CNN and multiple SPIDNA architectures). For each network the best run is surrounded by a square.

| input dimension | SNP encoding | Convolution type | Kernel size | Pooling size | Log-scaled output? | Sort chromosomes? | Use dropout? |
|---|---|---|---|---|---|---|---|
| $50 \times 400$ | 0/-1 | 1D | 2 | 2 | Yes | Yes | Yes |
| $50 \times 1784$ | 0/-1 | 1D | 2 | 2 | Yes | Yes | Yes |
| $50 \times 400$ | -1/1 | 1D | 2 | 2 | Yes | Yes | No |
| $50 \times 1784$ | -1/1 | 1D | 2 | 2 | Yes | Yes | No |

Table S1: Parameters used for the *Flagel* CNN

| | Method | Adaptive | Summary statistics | ABC correction | Alpha | Prediction error |
|---|---|---|---|---|---|---|
| 0 | ABC | No | No | No | No | 0.490 |
| 1 | ABC | No | No | Linear reg. | No | 0.357 |
| 2 | ABC | No | No | Ridge reg. | No | 0.363 |
| 3 | ABC | No | No | Single layer NN | No | 0.352 |
| 4 | MLP | No | No | No | No | 0.690 |
| 5 | MLP | No | Yes | No | No | 0.388 |
| 6 | Custom CNN | No | No | No | No | 0.485 |
| 7 | *Flagel* CNN 0/1 encoding | No | No | No | No | 0.447 |
| 8 | *Flagel* CNN 0/1 encoding | Downsampling | No | No | No | 0.437 |
| 9 | *Flagel* CNN -1/1 encoding | No | No | No | No | 0.719 |
| 10 | *Flagel* CNN -1/1 encoding | Downsampling | No | No | No | 0.984 |
| 11 | SPIDNA | No | No | No | No | 0.453 |
| 12 | SPIDNA | No | No | No | No | 0.637 |
| 13 | SPIDNA | Yes | No | No | No | 0.487 |
| 14 | SPIDNA | No | No | No | 0.849 | 0.592 |
| 15 | SPIDNA | Yes | No | No | 0.539 | 0.489 |
| 16 | ABC + SPIDNA | No | No | No | No | 0.462 |
| 17 | ABC + SPIDNA | No | No | Linear reg. | No | 0.364 |
| 18 | ABC + SPIDNA | No | No | Ridge reg. | No | 0.371 |
| 19 | ABC + SPIDNA | No | No | Single layer NN | No | 0.363 |
| 20 | ABC + SPIDNA | Yes | No | No | No | 0.458 |
| 21 | ABC + SPIDNA | Yes | No | Linear reg. | No | 0.363 |
| 22 | ABC + SPIDNA | Yes | No | Ridge reg. | No | 0.382 |
| 23 | ABC + SPIDNA | Yes | No | Single layer NN | No | 0.374 |
| 24 | ABC + SPIDNA | No | Yes | No | No | 0.476 |
| 25 | ABC + SPIDNA | No | Yes | Linear reg. | No | 0.339 |
| 26 | ABC + SPIDNA | No | Yes | Ridge reg. | No | 0.341 |
| 27 | ABC + SPIDNA | No | Yes | Single layer NN | No | 0.345 |
| 28 | ABC + SPIDNA | Yes | Yes | No | No | 0.474 |
| **29** | **ABC + SPIDNA** | **Yes** | **Yes** | **Linear reg.** | **No** | **0.335** |
| 30 | ABC + SPIDNA | Yes | Yes | Ridge reg. | No | 0.339 |
| 31 | ABC + SPIDNA | Yes | Yes | Single layer NN | No | 0.347 |

Table S2: Prediction errors of the best configuration of each method after hyperparameter optimization