

SPECIAL ISSUE

Automatic inference of demographic parameters using generative adversarial networks

Zhanpeng Wang¹ | Jiaping Wang¹ | Michael Kourakos² | Nhung Hoang² |
Hyong Hark Lee² | Iain Mathieson³ | Sara Mathieson¹ 

¹Department of Computer Science, Haverford College, Haverford, PA, USA

²Department of Computer Science, Swarthmore College, Swarthmore, PA, USA

³Department of Genetics, University of Pennsylvania, Philadelphia, PA, USA

Correspondence

Sara Mathieson, Department of Computer Science, Haverford College, Haverford, PA, USA.

Email: smathieson@haverford.edu

Funding information

National Institutes of Health, Grant/Award Number: R15HG011528 and R35GM133708

Abstract

Population genetics relies heavily on simulated data for validation, inference and intuition. In particular, since the evolutionary 'ground truth' for real data is always limited, simulated data are crucial for training supervised machine learning methods. Simulation software can accurately model evolutionary processes but requires many hand-selected input parameters. As a result, simulated data often fail to mirror the properties of real genetic data, which limits the scope of methods that rely on it. Here, we develop a novel approach to estimating parameters in population genetic models that automatically adapts to data from any population. Our method, *pg-gan*, is based on a generative adversarial network that gradually learns to generate realistic synthetic data. We demonstrate that our method is able to recover input parameters in a simulated isolation-with-migration model. We then apply our method to human data from the 1000 Genomes Project and show that we can accurately recapitulate the features of real data.

KEYWORDS

demographic inference, evolutionary modelling, generative adversarial network, simulated data

1 | INTRODUCTION

Simulation is a key component of population genetics. It helps to train our intuition and is important for the development, testing and comparison of inference methods. Because population genetic models such as the ancestral recombination and selection graphs (Griffiths & Marjoram, 1997; Neuhauser & Krone, 1997) are computationally intractable for inference but relatively easy to simulate, simulations are also heavily used for parameter inference. Approximate Bayesian Computation (ABC; Beaumont et al., 2002) is a widely used example. Regardless of the application, the goal is to simulate data that is 'realistic' in the sense that it resembles

real data from the population(s) of interest. Typically this is done by fixing some parameters that are fairly well-known, then choosing other parameters to match some property of the real data, usually based on summary statistics. However, this involves a potential loss of information in the reduction in summary statistics and then an implicit weighting on the relative importance of different summary statistics. Often, parameters that create simulations that match one type of summary statistic (e.g. the site frequency spectrum) do not match others (e.g. linkage disequilibrium patterns; Beichman et al., 2017). Here, we present a novel parameter learning approach using Generative Adversarial Networks (GANs). Our approach creates both realistic simulated data and a quantitative way of determining

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2021 The Authors. *Molecular Ecology Resources* published by John Wiley & Sons Ltd.

the match between any simulations proposed for a particular real data set. For us, 'realistic' means 'cannot be distinguished from real data by a machine learning algorithm', specifically a convolutional neural network (CNN).

Machine learning (ML) methods have been emerging more broadly as promising frameworks for population genetic inference. The high-level goal of training a ML method is to learn a function from the input (genetic data) to the output (evolutionary parameters). Some early efforts used machine learning to account for issues that arise with high-dimensional summary statistics (Blum & François, 2010; Sheehan & Song, 2016; Ronen et al., 2013). More recently, machine learning approaches have used various forms of convolutional, recurrent and 'deep' neural networks to improve inference and visualization (Adrion et al., 2020; Battey et al., 2021; Gower et al., 2020; Flagel et al., 2019; Sanchez et al., 2020; Torada et al., 2019; Chan et al., 2018). One of the goals of moving to these approaches was to enable inference frameworks to operate on the 'raw' data (genotype matrices), which avoids the loss of information that comes from reducing genotypes to summary statistics. However, these algorithms rely heavily on simulated data sets for training. In machine learning more broadly, data are often hand-labelled with 'true' values—part of these data are used to train the model, and part are held aside to test the model. In population genetics, such 'labelled' training data are extremely limited, because the evolutionary ground truth is rarely known with certainty. Thus, all approaches rely on simulations to train and validate ML models.

Current simulators (Kern & Schrider, 2016; Excoffier et al., 2013; Hudson, 2002; Haller & Messer, 2019; Kelleher et al., 2016; Ewing & Hermisson, 2010; Teshima & Innan, 2009) are well equipped to replicate mechanisms of evolution but require many user-selected input parameters including mutation rates, recombination and gene conversion rates, population size changes, natural selection, migration rates and admixture proportions. We do not always have a good sense of what these parameters should be, especially in understudied populations and non-model species. For example, mutation and recombination rates estimated in one population are frequently used to simulate data for another, despite the fact that these rates differ between populations (Adrion, Cole, et al., 2020; Harris, 2015; Hinch et al., 2011; Harris & Pritchard, 2017; Kessler et al., 2020).

Generative models provide one route to simulating more realistic population genetic data. Typically, generative models create artificial data based directly on observed data, without an explicit underlying model. They have been used to create synthetic examples in a wide range of fields, from images and natural language to mutational effects (Riesselman et al., 2018) and single cell sequencing (Lopez et al., 2018). In particular, Generative Adversarial Networks (GANs) work by creating two networks that are trained together (Mirza & Osindero, 2014; Goodfellow et al., 2014). One network (the *generator*) generates simulated data, while the other network (the *discriminator*) attempts to distinguish between 'real' data and 'fake' (synthetic) simulations. As training progresses, the generator learns more about the real data and gets better at creating realistic examples, while the discriminator learns to pick up on subtle differences

and gets better at distinguishing examples. After training is complete, the generator can be used to create new examples that are indistinguishable (by the discriminator) from real data, but where the ground truth is known (i.e. labelled data).

The use of GANs in population genetics is just beginning. Recently, Yelmen et al. (2021) created a GAN that generates artificial genomes that mirror the properties of real genomes. Their approach does not include an evolutionary model, so the resulting artificial genomes are 'unlabelled'. Such an approach is useful for creating proxy genomes that preserve privacy but still maintain realistic aggregate properties. However, this synthetic data could not be used downstream to train or validate supervised machine learning methods since no evolutionary ground truth is known.

Here, we present a parametric GAN framework that combines the ability to create realistic data with the interpretability that comes from an explicit model of evolution. The discriminator is a permutation-invariant CNN that takes as input a genotype matrix (representing a genomic region) and classifies it as real data or synthetic data. Throughout training, the discriminator tries to get better at this binary classification task. The generator is a coalescent simulator that generates genotype data from a parameterized demographic history. The generator is trained using a simulated annealing algorithm that proposes parameter updates leading to more discriminator confusion. The discriminator is trained using a gradient descent approach that is standard for neural networks. We apply our method, called *pg-gan*, in a variety of scenarios to demonstrate that it is able to recapitulate the features of real genetic data and confuse a trained discriminator. Although we focus on humans, the underlying methodology enables the simulation of any population or species, regardless of how much is known *a priori* about their specific evolutionary parameters.

We anticipate that the approach outlined in this work will be useful in strengthening the match between simulated and real data, especially for understudied populations that deviate from broad geographic groups. In addition, our discriminator can be used on its own (after training) to evaluate and compare different candidate simulations for the same real data set. Downstream, our simulations can be used as a starting point for other methods that seek to quantify local evolutionary forces such as natural selection or mutation rate heterogeneity. There has also been a push in the population genetics community to standardize simulation resources (Adrion et al., 2020)—we see our method as contributing to the assessment and refinement of published models as they are applied to new data sets.

2 | MATERIALS AND METHODS

At a high level, our method works by simulating data from an underlying evolutionary model, then comparing it to real data via a neural network discriminator. As the discriminator is trained, it tries to minimize a loss function that incentivizes learning the difference between real data and synthetic data. But at the same time, the

generator refines the evolutionary model so that it recapitulates the features of real data and attempts to confuse the discriminator. At the end, the evolutionary model can be used to simulate additional realistic data for use in downstream applications or method comparisons. Additionally, the final parameters of the evolutionary model can be interpreted to learn more about the population or species of interest.

A GAN is not a traditional optimization problem—due to the dual nature of the generator and discriminator there are two optimization problems in a minimax framework, and it is difficult to evaluate the final trained model. Often the ‘GAN confusion’ (discriminator classification accuracy) can be used to assess the success of the algorithm—a high classification accuracy (close to 1) indicates that the simulations are not capturing the real data and the discriminator is easily able to tell the difference between the two types of data. A low classification accuracy (close to 0.5) ideally indicates the evolutionary model has created simulations that are well-matched to the real data. However, an accuracy close to 0.5 could also mean that the discriminator has not learned anything and is either flipping a coin when classifying examples, or classifying all examples as the same class.

Training a GAN is a delicate balance. If the discriminator learns too quickly and becomes very good at identifying a specific setting of the simulated data from the real, then all proposals by the generator may look equally confusing. As a result, many generator proposals will be rejected and the discriminator will simply keep getting better at distinguishing the current setting from real data. On the other hand, if the discriminator learns too slowly, it may not be able to identify any generator proposals as better or worse. This often leads to a ‘random walk’ across the parameter space, with the discriminator classifying everything as real or everything as simulated regardless of the generator’s proposals. Throughout the Methods section, we outline techniques and strategies for balancing training and identifying degenerate states.

In the Method subsections below, we first outline the notation for *pg-gan* and discuss the general training strategy. Then, we provide further details about the generator and discriminator architectures. Finally, we discuss applications of *pg-gan* to both simulated and real training data, as well as methods for evaluating the performance.

There are two inputs to the method, shown in orange in Figure 1. The first input is an evolutionary model parameterized by vector Θ . The parameters can be very flexible, including evolutionary event times, effective population sizes and rates of mutation, recombination, migration and exponential growth. The parameters Θ can be fed into the generator G to produce a simulated region z , which we write as

$$z \sim G(\Theta).$$

The second input is a set of real data. We use x to denote a generic region from the real data. Both z and x have the same shape $(n, S, 2)$ where n is the number of haplotypes, S is the number of SNPs

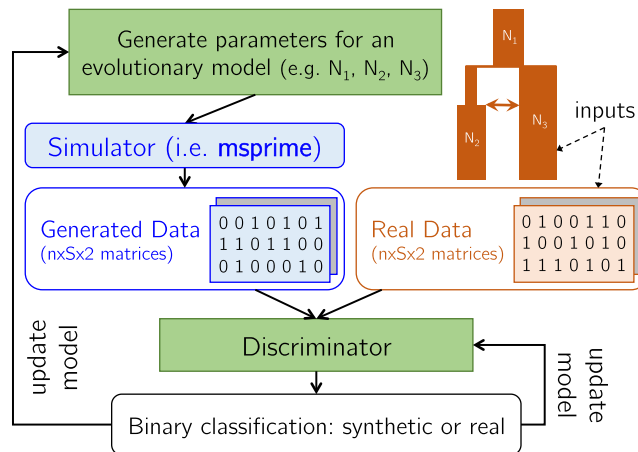


FIGURE 1 *pg-gan* algorithm overview. The inputs to our method are an evolutionary model and a set of real data (orange). The parameters of the generator and discriminator (green) are updated in a unified training framework using simulated annealing (generator) and backpropagation (discriminator). The generated data and real data are analysed one genotype matrix at a time, where n is the number of haplotypes and S is the number of SNPs retained in each region. Inter-SNP distances are also fed in as a second channel, which provides the discriminator with information about SNP density

in the region. The first channel represents the genotypes, and the second channel represents the inter-SNP distances. The outputs of *pg-gan* are the optimal evolutionary parameters Θ^* for the generator G , and a binary classifier D (the discriminator) which can predict if genomic regions are real or fake. Specifically, $D(x)$ is the predicted probability that region x is *real*.

To incentivize the competing goals of the generator and discriminator, we minimize binary cross-entropy loss functions. If we have M regions of simulated data $\{z^{(1)}, \dots, z^{(M)}\}$ generated under $G(\Theta)$, then the generator loss function is

$$\mathcal{L}_G(\Theta) = -\frac{1}{M} \sum_{m=1}^M \log D(z^{(m)}).$$

This loss function is cross-entropy, but where we only have one class (the generated data), which we want the discriminator to classify as real (label 1).

At the same time, the discriminator D is trying to classify the generated regions as fake (label 0) and the real regions as label 1. Therefore, the discriminator loss function for M regions of real data $X = \{x^{(1)}, \dots, x^{(M)}\}$ and M regions of simulated data $\{z^{(1)}, \dots, z^{(M)}\}$ generated under $G(\Theta)$ is

$$\mathcal{L}_D(\Theta, X) = -\frac{1}{M} \sum_{m=1}^M \left[\log D(x^{(m)}) + \log(1 - D(z^{(m)})) \right].$$

Algorithm 1. (in the style of (Goodfellow et al., 2014)) shows the overall training of *pg-gan*.

Algorithm 1: Training pg-gan.

Input: evolutionary model parameterized by Θ , real data X , architectures of generator G and discriminator D
Output: optimal parameters Θ^* (which can be fed into G to produce synthetic data), trained discriminator D

for num pre-training iterations **do**

 Initialize evolutionary parameters randomly to Θ
 Train discriminator using real data and simulated data under $G(\Theta)$

end

Initialize $\Theta^{(0)}$ to the parameters that caused the highest accuracy during pre-training

for each training iteration i **do**

 Use simulated annealing to select parameters Θ near $\Theta^{(i)}$ that minimize the generator loss

$$\mathcal{L}_G(\Theta) = -\frac{1}{M} \sum_{m=1}^M \log D(z^{(m)}),$$

 where $\{z^{(1)}, \dots, z^{(M)}\}$ is a mini-batch of data simulated under $G(\Theta)$

if best parameters Θ are accepted **then**

for num mini-batches **do**

 Sample mini-batch of M regions $X = \{x^{(1)}, \dots, x^{(M)}\}$ from the real data

 Simulate mini-batch of M regions $\{z^{(1)}, \dots, z^{(M)}\}$ under $G(\Theta)$

 Update the discriminator by minimizing binary cross-entropy loss (below) using gradient descent

$$\mathcal{L}_D(\Theta, X) = -\frac{1}{M} \sum_{m=1}^M [\log D(x^{(m)}) + \log(1 - D(z^{(m)}))]$$

end

 Set $\Theta^{(i+1)} \leftarrow \Theta$

end

end

Run the above procedure K times; select the parameters that minimize discriminator accuracy on the simulations.

2.1 | Generator

In image and video generation, the generator often takes the form of a CNN, since a large array of pixel information must be generated from a low-dimensional vector of noise (see Figure 1 of (Radford et al., 2015) for the architecture of a CNN-based image generator). For our purposes, we do not need to generate the individual genotypes for each training example, but we do need to generate candidate parameters for input into an evolutionary simulator (we use `msprime` (Kelleher et al., 2016) in this study).

Using this lens, we can view the generator learning problem as minimizing the multivariate generator loss function $\mathcal{L}_G(\Theta)$ with respect to Θ . We optimize the loss using simulated annealing (Pincus, 1970) due to its flexible parameter updates and lack of reliance on an analytic gradient. In simulated annealing, initial parameter values are proposed and then gradually refined. A *temperature* is used to control whether or not new parameter proposals are accepted. The temperature usually begins at a high value, indicating that sub-optimal parameter choices may be accepted liberally to facilitate exploration of the parameter space. As training proceeds, the temperature ‘cools’, reducing the chance of accepting a poor parameter choice and allowing the method to converge on a set of parameters that optimizes the desired function. Unlike ABC methods which require simulating from the entire parameter space before analysing the real data, this simulated annealing approach uses the real data to adaptively narrow the focus to promising regions of the search space.

We use a pre-training phase (described in the *Discriminator* subsection) to choose a starting value for each evolutionary parameter, which forms the initial parameter vector $\Theta^{(0)}$. We set the temperature for simulated annealing $T^{(0)} = 1$ and linearly decrease it to 0

over a fixed number of iterations. During each training iteration i , several new sets of candidate parameters are proposed, and evaluated based on the generator loss function $\mathcal{L}_G(\Theta)$. Each new set of parameters is proposed by sampling from a normal distribution around each current value, with variance based on the temperature. This allows the algorithm to explore the parameter space quickly in the beginning and refine the estimates towards the end of GAN training. More formally, at iteration i , the candidate proposal for parameter p is

$$\Theta_p^{(\text{proposal})} \sim \mathcal{N}(\Theta_p^{(i)}, \sigma_p^2 \cdot T^{(i)})$$

where σ_p^2 is the initial variance, which is based on the range of plausible values for each parameter. Out of the several candidate proposals, we choose the one that minimizes $\mathcal{L}_G(\Theta)$. Then, we compare this loss to the loss of the previous iteration. If the proposal reduces or maintains the generator loss, we always accept it. If not, we use the simulated annealing temperature to help define a threshold for acceptance. Formally, if the proposal is Θ and the current set of parameter values at iteration i is $\Theta^{(i)}$, then the acceptance probability is

$$p_{\text{accept}} = \frac{\mathcal{L}_G(\Theta^{(i)})}{\mathcal{L}_G(\Theta)} \cdot T^{(i)}.$$

If the proposed parameters are accepted we train the discriminator using several mini-batches (with the simulated regions generated under Θ), then we set $\Theta^{(i+1)} \leftarrow \Theta$. An important point is that we do not train the discriminator using the new parameter proposal unless it is accepted. During the candidate proposal phase, we are evaluating the parameter choices through the generator loss only.

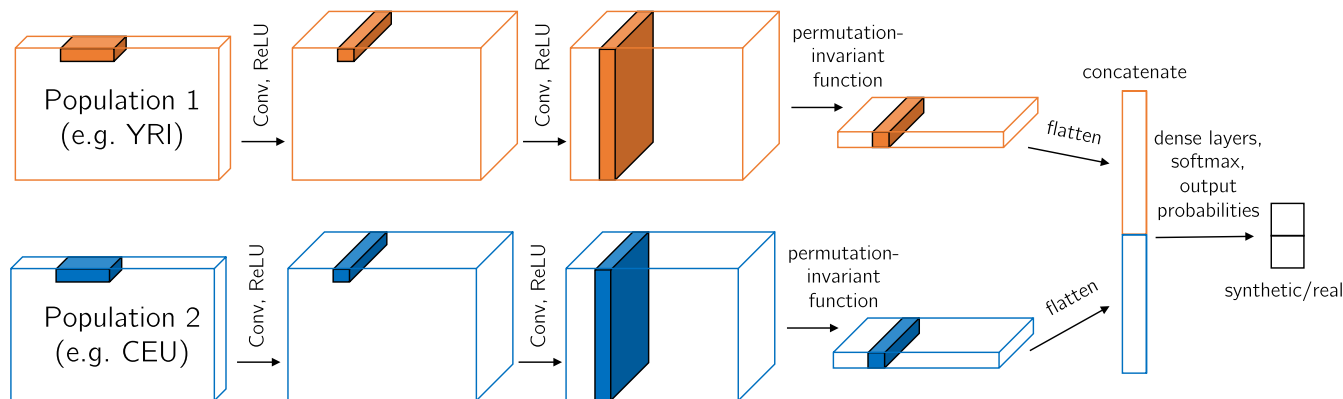


FIGURE 2 Multi-population CNN discriminator architecture. Each example region is of shape $(n, S, 2)$, where n is the number of haplotypes (usually with $n/2$ from population 1 and $n/2$ from population 2). The convolutional filters for population 1 and 2 are shared (i.e. not separate weights) so that haplotype commonalities can be more easily identified. The final output of the discriminator is the probability the region is real (which can be subtracted from 1 to find the probability the region is simulated). This CNN can be reduced for one population or extended for three populations

2.2 | Discriminator

For the architecture of the discriminator, we use a permutation-invariant CNN based on `defiNETti` (Chan et al., 2018). Each region x (real) or z (simulated) has shape $(n, S, 2)$ where n is the number of haplotypes in the sample, S is the number of retained SNPs, and 2 indicates there is one channel for the genotypes and one channel for inter-SNP distances. The inter-SNP distances are duplicated down each column to allow this slice of the tensor to have the same shape as the genotype information. This also ensures that each convolutional filter processes the genotypes and associated distances at the same time. Alternatively, the convolutional layers can be used on the genotypes only, and the distances concatenated later as a vector. However, this approach does not allow the processing of the two channels to be as tightly coupled. We use convolutional filters of shape 1×5 (1 haplotype, 5 SNPs) to ensure that the order of haplotypes does not impact the results. We use ReLU as the activation function for all layers and also use dropout (Srivastava et al., 2014) during training to guard against overfitting. After several convolutional layers, we condense the output by applying a column-wise permutation-invariant function. We experimented with both `max` and `sum` as permutation-invariant functions and decided to use `sum` throughout. It generally causes the discriminator to learn more slowly than `max`, allowing the generator time to find good parameter choices. `max` sometimes causes the discriminator to converge quickly, easily distinguishing the real from simulated data before the generator can move to a promising location in the parameter space. Note that we need a fixed number of SNPs in each region to make sure the discriminator output is always of the same size. However, we do not need a consistent number of haplotypes, provided that the permutation-invariant function used is not sensitive to this number (i.e. `max` or `avg` would be fine but `sum` would not).

For models that consider multiple populations, we augment this framework to include separate permutation-invariant components for each population, then concatenate the flattened output before

input into the dense layers at the end of the network. An illustration of our discriminator architecture for two populations is shown in Figure 2.

Through discriminator training, we seek to minimize the loss function $\mathcal{L}_D(\Theta, X)$, with a small entropy term subtracted to disincentivize predicting all the same class. This entropy term is different from the entropy regularization used to prevent mode collapse (Dieng et al., 2019), a common problem in GAN training. In such cases, the goal is to increase the entropy of the *generator* so it can produce a multi-modal distribution (e.g. different types of images such as hand-written digits). Mode collapse is not an issue for `pg-gan`, as a single set of evolutionary parameters is desired. To minimize our discriminator loss function, we use gradient descent (via backpropagation) with mini-batches of 100 training examples (half are real and half are simulated). For each training iteration, we perform 100 mini-batch training updates if the proposed parameters are accepted. This allows the discriminator to learn gradually, as the parameters are being refined. While a classification accuracy close to 0.5 is desired by the end of training, the discriminator accuracy may be close to this value early on in the training process simply because it has not learned anything yet. The goal is for the discriminator to be optimized to distinguish real from simulated data as much as possible and *still* be wrong half the time.

Due to the simulated annealing training of the generator, initial step sizes of the parameters can be large to explore the parameter space more quickly. This can present a problem for discriminator training. If the parameters change too quickly, the discriminator does not have time to learn the difference between the real data and data simulated under a wide variety of parameter setting. In some situations, this leads the discriminator to fail to learn anything and it predicts the same class (either real or fake) for all regions. To combat this issue, before GAN training begins we pre-train the discriminator only, using a variety of randomly sampled parameter values. We find that pre-training gives the discriminator an overall sense of the data, so that when generator training begins the discriminator is able to

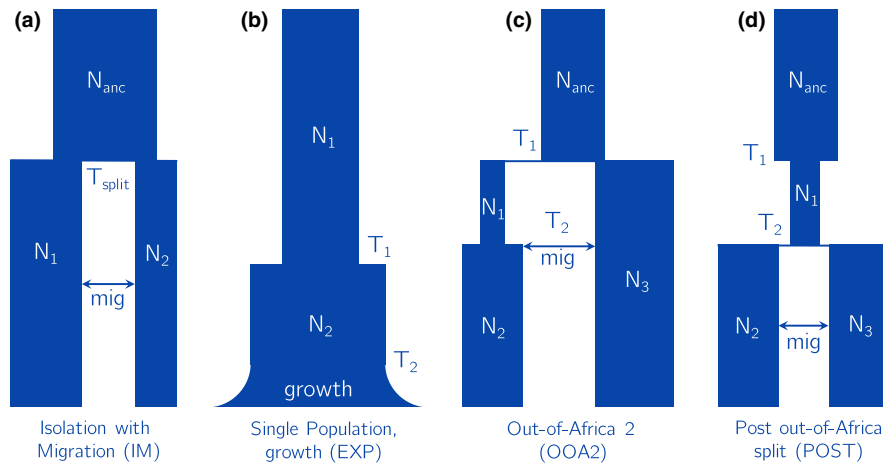


FIGURE 3 Set of models. (a) A six-parameter, two-population isolation-with-migration model, which we use in the simulation study. The migration event is a single pulse at time $T_{\text{split}}/2$, and can be in either direction. The final parameter (not shown in this diagram) is the recombination rate. (b) A five-parameter, single-population exponential growth model, which we use to infer histories for YRI, CEU and CHB separately. (c) A seven-parameter, two-population model, which we fit separately for YRI/CHB and YRI/CEU. The migration can be in either direction. (d) A seven-parameter, two population model which we fit to CEU/CHB. Migration occurs at $T_2/2$ and can be in either direction

identify which generated regions were closer to the real and which were further away. We run the combined (generator and discriminator) training for 300 iterations.

2.3 | Simulation study

To validate our approach, we first select the training data set to be a simulated one, so that we can test whether the inferred parameters are correct. To assess a variety of different types of parameters, we choose an isolation-with-migration model (see Figure 3a) with six parameters. The parameters include three effective population sizes: N_{anc} for the ancestral population size, and N_1 and N_2 for the sizes of each population after the split. We also infer the split time T_{split} , and the strength of a directional migration pulse (mig) at time $T_{\text{split}}/2$. Finally, we infer the per-base, per-generation recombination rate (reco). We evaluate the inferred parameters based on how well they match the true parameters. See Table 1 for the ranges and units of each parameter.

2.4 | 1000 Genomes data analysis

To demonstrate the effectiveness of our method on real data, we use *pg-gan* to infer demographic parameters for both single- and multi-population models in humans. To ensure that the real data are as similar as possible to the simulated data, we run several pre-processing steps. To avoid processing the real data on-the-fly during training, we follow a data extraction pipeline to convert the real data into HDF5 format (Miles, 2015, 2017). Before converting VCF information into HDF5 format, we select haplotypes from each population and filter non-segregating and multi-allelic sites. The number of haplotypes is flexible (due to the permutation-invariant framework).

TABLE 1 Parameter ranges

Parameter	Min	Max	Units
N_e	1000	30,000	Individuals
reco	1×10^{-9}	1×10^{-7}	Per base per generation
mut	1×10^{-9}	1×10^{-7}	Per base per generation
N_{anc}	1000	25,000	Individuals
T_{split}	500	20,000	Generations
mig	-0.2	0.2	Fraction of individuals
N_1	1000	30,000	Individuals
N_2	1000	30,000	Individuals
growth	0	0.05	Per generation
N_3	1000	30,000	Individuals
T_1	1500	5000	Generations
T_2	100	1500	Generations

When inferring a parameter, we initialize its value by drawing a value uniformly from the given ranges. For each parameter update, we do not allow the parameter to go up to or outside its range. Overall, the ranges are meant to be plausible values based on previous studies or reasonable evolutionary events.

We use between 196 and 198, matching the minimum number of individuals in each 1000 Genomes population.

During training, for each region x fed into *pg-gan*, we select a start SNP randomly from the entire genome. This random start point mitigates the effects of correlated nearby regions and local variations in mutation and recombination rate. Starting with this SNP, $S = 36$ biallelic SNPs are retained (along with their inter-SNP distances), which means the region has a flexible length. If 36 SNPs would cause the region to extend past the end of a chromosome,

we reject the start SNP and sample a new one. For each region, we retain it if at least 50% of the bases are inside callable regions, as defined by the '20120824' strict mask (1000 Genomes Project Consortium, 2015).

For both the real and simulated data, we recode the genotypes by setting the minor allele to the value '1' and the major allele to the value '-' so that the discriminator cannot learn to distinguish real data based on reference bias or ancestrally misidentified states. For the simulations where we must specify a region length L , we choose $L = 50\text{kb}$, which ensures that in the majority of situations we have at least $S = 36$ SNPs. The middle 36 SNPs are retained, and any regions with insufficient SNPs are centred and zero-padded. Such regions would automatically look very different from the real data, so the generator quickly learns to avoid parameters that cause insufficient SNPs.

We test four models: Figure 3b–d, as well as a three-population Out-of-Africa model originally specified in (Gutenkunst et al., 2009). The single-population model has five parameters: two effective population sizes N_1 and N_2 , two size-change points T_1 and T_2 , and the rate of exponential growth in the recent past. We fit this model to three human populations from the 1000 Genomes project: YRI (West African), CEU (European), and CHB (East Asian). The second model (OOA2) is a simplified two-population Out-of-Africa model. There are seven parameters: four effective population sizes, two time-change points and a migration pulse that can be in either direction, allowing for migration between African and non-African populations. We fit this model to two pairs of populations: YRI/CEU and YRI/CHB. The third model (POST) represents the post-out-of-Africa split between the ancestors of Europeans and East Asians. In this seven-parameter model, we allow a pre-split bottleneck and directional migration. We fit this model to the pair of populations CEU/CHB. Finally, we apply the three-population Out-of-Africa model (OOA3) to YRI/CEU/CHB, as implemented in `stdpopsim` (Adrion et al., 2020).

2.5 | Evaluation metrics

One pervasive issue with GANs is the lack of a natural evaluation metric (see (Borji, 2019) for a comprehensive overview of GAN evaluation metrics). Many GANs have been evaluated qualitatively through user studies designed to see whether humans find the generated data realistic (Xu et al., 2018). For images, videos or text, this type of evaluation can be informative (although it tends to favour generators that memorize specific real examples; Borji, 2019), but this is not directly possible in the case of genetic data.

Visualizing summary statistics is an alternative, although since we do not know which statistics are sufficient for the model, it is dangerous to rely on these alone as a final evaluation metric. It is possible the discriminator is learning other statistics or representations of the data that we are not aware of. In addition, explicitly matching some types of statistics can bias the resulting fitted model. For example, Beichman et al. (2017) found that SFS-matching methods like `daði` (Gutenkunst et al., 2009) and `SMC++`

(Terhorst et al., 2017) are not able to recapitulate LD statistics. Further, we currently do not have an exhaustive or sufficient set of summary statistics that could be used to identify model parameters directly in a likelihood framework. However, as a qualitative assessment of our results, we compare summary statistics computed on the real data and data simulated under our inferred parameters. This gives us a sense of which features of real data agree with our simulations and which do not.

To that end, we use seven types of summary statistics. In all cases, we use 5000 regions of real data (chosen randomly) and 5000 regions of simulated data (each simulated independently under our inferred parameters) to compute the statistics. All pre-processing is the same as for GAN training, except for Tajima's D where we fix the region length, not the number of SNPs.

- **SFS:** We compute the site frequency spectrum (SFS) by counting the number of singletons, doubletons, etc in each of 5000 regions of real and simulated data. We plot the first 10 entries.
- **Inter-SNP distances:** We plot the distribution of inter-SNP distances for both the real and simulated data (measured in base pairs). This provides a general measure of SNP density.
- **LD:** We compute linkage disequilibrium (LD) by clustering pairs of SNPs based on their inter-SNP distance. We divide these distances into 15 bins and average the correlation r^2 within each one.
- **Pairwise heterozygosity:** We plot the distribution of pairwise heterozygosity (π), computed separately for each region.
- **Tajima's D :** We plot the distribution of Tajima's D , computed separately for each region. Here, we fix the region length to $L = 50\text{ kb}$ instead of fixing the number of SNPs, as otherwise the distribution would be the same as pairwise heterozygosity.
- **Number of haplotypes:** We plot the distribution of number of haplotypes for each region.
- **Hudson's F_{st} :** For the two-population split models, we use F_{st} to measure population differentiation (Hudson et al., 1992).

As a more quantitative evaluation, we also report the final discriminator classification accuracy. However, even this metric is not easy to interpret, as an accuracy close to 0.5 may indicate a degenerate situation where the discriminator has not learned anything (see Figure S1 for an example). Thus, for each model and data set we run `pg-gan` $K = 10$ times and select the model that minimizes the classification accuracy of the discriminator on the final *generated* data (not using any real data). The more generated data that the discriminator classifies as real—that is, the lower the discriminator accuracy—the better the generator. This metric was inspired by the Inception Score (Salimans et al., 2016) used to evaluate GANs, where generated data are fed into a more powerful discriminator. Since no generated region has ever been seen by the discriminator before, all generated regions are implicitly 'test data'. In this way, we avoid relying on a held-aside real data set for evaluation, which allows us to use the (limited) real data exclusively for training.

Finally, we also ran a comparison study against the ABC method `fastsimcoal` (Excoffier & Foll, 2011). We provided `fastsimcoal`

with three of the same models (IM, OOA2, and POST) as well as a simulated joint SFS, and the full genome joint SFS for YRI/CEU, YRI/CHB, and CEU/CHB. Then, we compared the parameters and summary statistics from *fastsimcoal* to those from *pg-gan*.

3 | RESULTS

3.1 | Simulation study

To validate our method, we first simulated the training data set, so we knew the true evolutionary parameters. We fit the six-parameter IM model from Figure 3a, using the parameter ranges in Table 1. Throughout, we usually fix the mutation rate to 1.25×10^{-8} per base per generation, but it could be inferred along with the other parameters in species or populations where it is less established. See Figure S2 for example where we infer mutation rate as well as the other six parameters of the IM model.

During the pre-training phase, we train the discriminator on up to 10 different parameter sets, randomly chosen from the ranges in

Table 1. We select $\Theta^{(0)}$ to be the first set that achieves at least 90% discriminator classification accuracy (or the set that maximizes accuracy in the case when we do not achieve 90% after 10 pre-training iterations). This enables the discriminator to gain some structure that is relevant to the data before combined training begins. During each main training iteration, we choose 10 independent proposals for each parameter, keeping the other parameters fixed. This creates $10 \times P$ possible parameter sets, where P is the number of parameters ($P = 6$ for the IM model). We select the set that minimizes the generator loss, which has the effect of modifying one parameter each iteration. We also tested modifying all the parameters each iteration, but generally found that updating one at a time led to more stable and consistent results. For each parameter p , we set the initial variance σ_p^2 to the parameter range divided by 15.

We performed 10 independent initializations of *pg-gan* on the full set of six parameters for the IM model. We selected the results that minimized discriminator accuracy on the final generated data. The results are shown in Figures 4,5 and Table 2. The first subplot in Figure 4 shows the losses for both the generator and discriminator. Since the generator loss considers half as many regions, it is multiplied

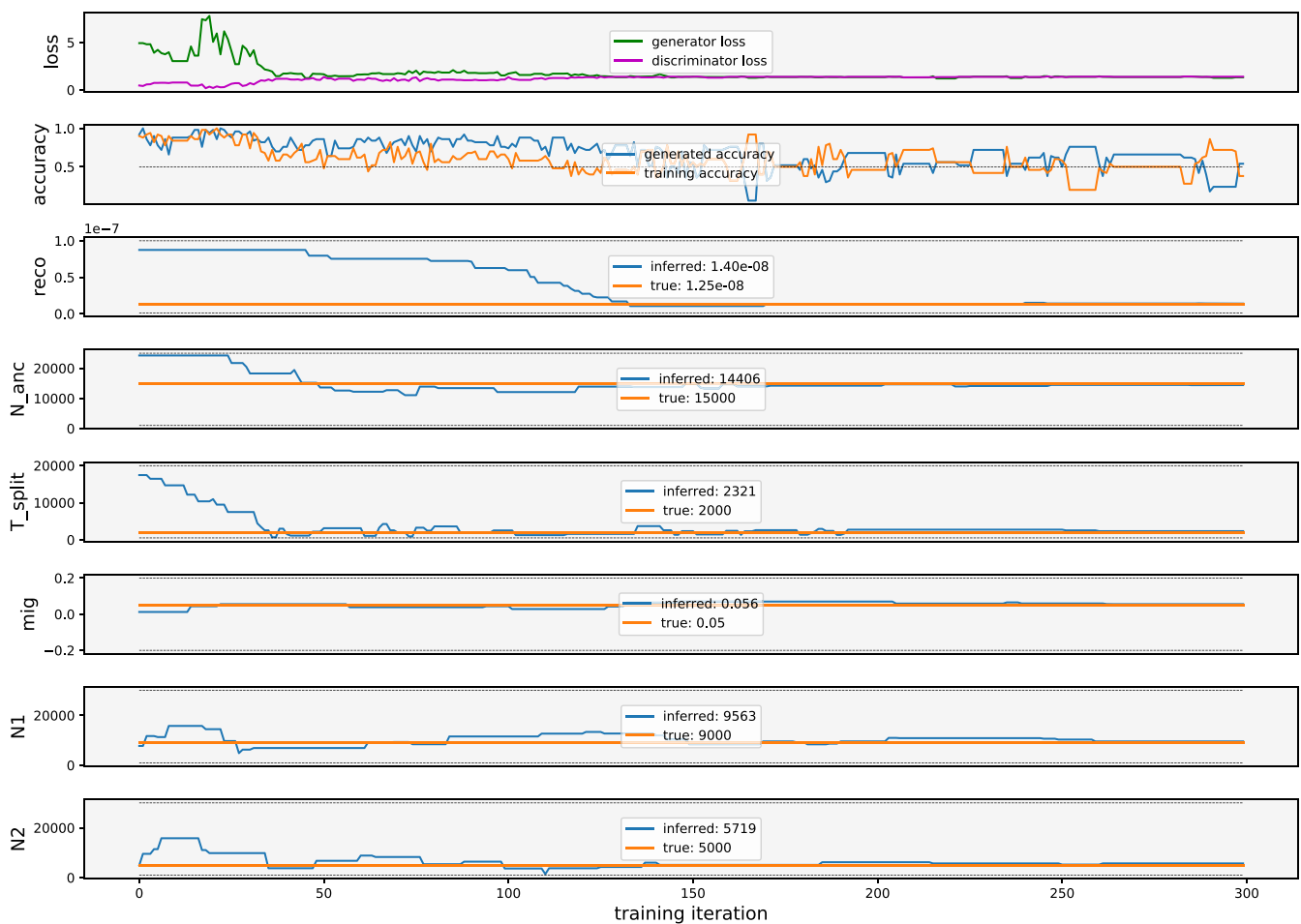


FIGURE 4 IM model parameter inference on simulated training data. In this scenario, we jointly infer the six parameters of the IM model from Figure 3a. The top plot shows both loss functions over the course of GAN training, and the second plot shows classification accuracy for both simulated and training data. The remaining plots show the model parameters as they are refined throughout GAN training. The inferred values are taken at the final iteration

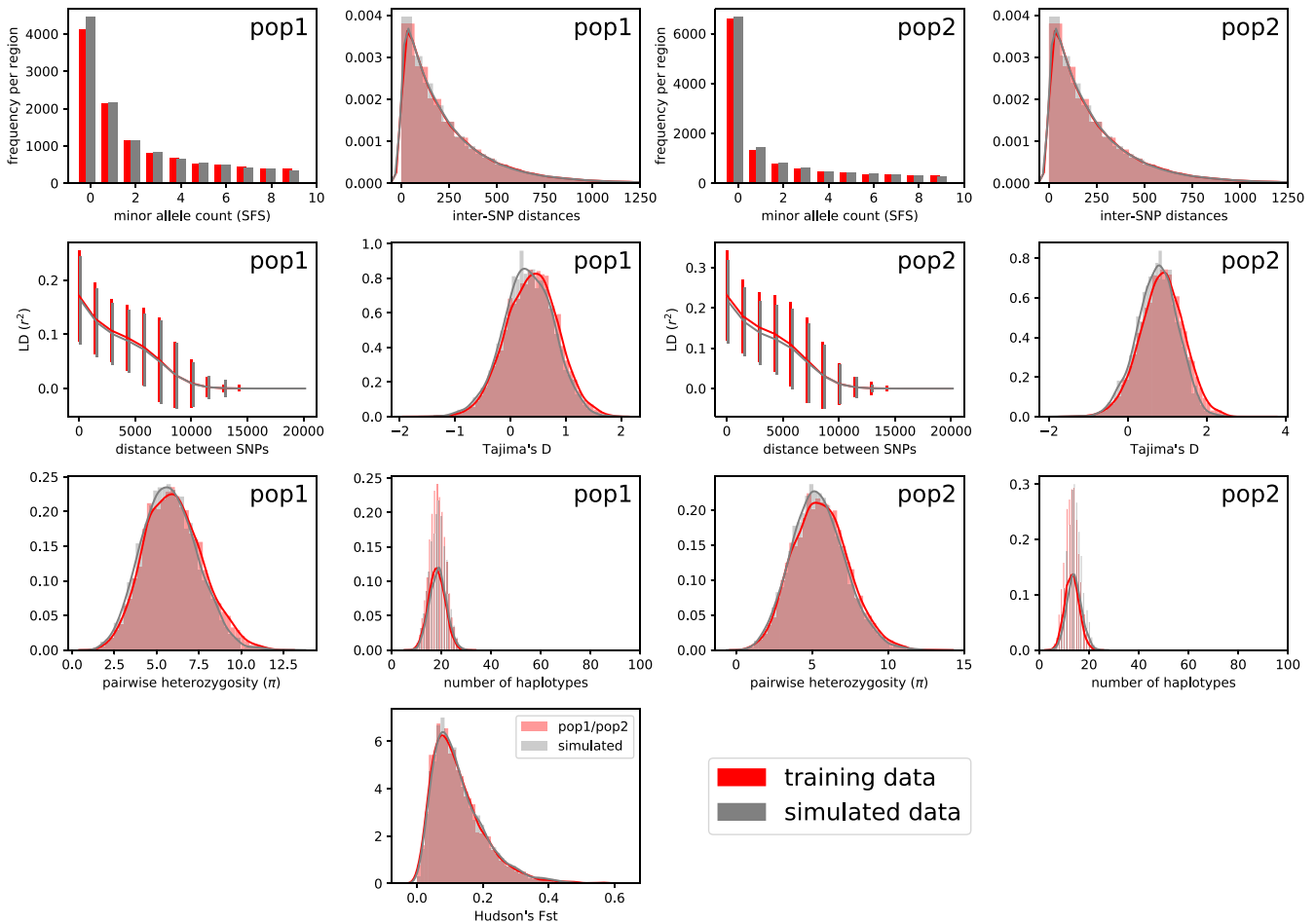


FIGURE 5 IM model statistics on simulated training data. Summary statistics for data simulated under our inferred parameters ('simulated data'), compared with data simulated under the true parameters ('training data'). Subfigures on the left correspond to statistics from the first population, and those on the right correspond to the second population. In the bottom panel, we show F_{st} between the two populations

TABLE 2 Comparison of *pg-gan* and *fastsimcoal*

	N_1	N_2	N_{anc}	T_{split}	mig	reco
TRUE	9,000	5,000	15,000	2,000	0.050	1.25e-08
<i>pg-gan</i>	9,563	5,719	14,406	2,321	0.056	1.40e-08
<i>fastsimcoal</i>	8,455	4,864	15,395	1,887	0.028	-

Inferred parameters for the IM model (see Figure 3a). *pg-gan* results correspond to Figures 4,5, and *fastsimcoal* results correspond to Figure 6. For *fastsimcoal*, recombination rate is not shown, since it cannot be inferred from the SFS.

by two to be on same scale as the discriminator loss. At first, the generator loss is high and the discriminator loss is low, because the discriminator is easily able to detect the difference between simulated and training data. The second plot shows the discriminator accuracy on both simulated and training data. Both accuracies are initially high and then reduce to around 0.5. We see that here, *pg-gan* is able to find parameter values that bring the discriminator close to an accuracy of 0.5. The final classification accuracy on generated data here was 0.54, and the overall accuracy (considering both generated and training data) was 0.46. Our inferred parameter values are close to

the true values (Table 2) and the site frequency spectrum and other summary statistics of data simulated with these parameter values closely match the summary statistics of the training data (Figure 5).

As a comparison, we performed ABC inference using *fastsimcoal* using the same IM model, fitting the joint SFS from data simulated under the true model parameters. *fastsimcoal* closely matches the SFS and true parameters (Figure 6, Table 2), although it is not able to infer recombination rate. If we give it the correct recombination rate, then it closely matches the other summary statistics in Figure 5 as well.

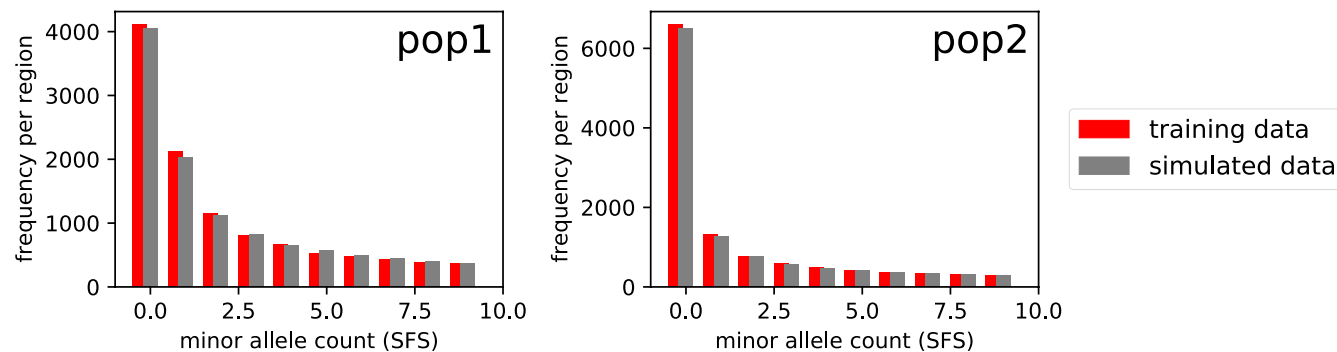


FIGURE 6 IM model SFS as inferred by *fastsimcoal*. Here, we compare the true SFS ('training data') with the SFS computed from data simulated under the parameters learned by *fastsimcoal* ('simulated data')

3.2 | 1000 Genomes data analysis

We analysed three populations (YRI, CEU and CHB) separately, each under the five-parameter model with recent exponential growth (EXP; Figure 3b). For all single-population results, we used $n = 198$ (size of CEU) and $S = 36$. Unlike the simulated example above, we fix the distribution of recombination rates by sampling from the distribution of HapMap combined recombination rates (International HapMap Consortium, 2007), whereas in principle this distribution could also be inferred.

To assess the impact of the model on our evaluation metrics (classification accuracy and summary statistics), we first fit a one-parameter demographic model with a single constant population size N_e . We then contrast this result with the five-parameter exponential growth model (EXP). The summary statistics for these results are shown in Figure 7 for YRI and CHB, and a summary for all populations is shown in Figure 8a. Inferred parameters for each population under the five-parameter exponential growth model are shown in Table 3. The effect of the Out-of-Africa bottleneck (N_2) is very apparent in CEU and CHB, but absent in YRI. Data simulated with fitted parameters for YRI contain many more singletons than the real data, possibly indicating the recent exponential growth rate (or the time of onset T_2) is overestimated. On the other hand, low power to detect rare variants in the real data could explain a lack of singletons in YRI or other populations.

We also compared summary statistics (Methods) between the real data and data simulated under the parameter choices corresponding to the two scenarios from Figure 8a. In Figure 7, we show two sets of summary statistics each for YRI and CHB. On the left, we show the one-parameter demography results, and on the right we show the five-parameter results (using HapMap recombination rates in both cases). While some statistics match closely, others are less well-matched, consistent with the discriminator being imperfectly confused. Summary statistics for CEU are shown in Figure S3. For CEU, both the one- and five-parameter models produced low classification accuracy, but the summary statistics are imperfect. This likely indicates that the discriminator did not learn as well in this scenario, not that the generator is producing high-quality simulated data.

For all our results, we discard any run where the discriminator classifies all regions in the same way (either all real or all simulated)

at the end of training. For each set of 10 runs, 0–2 runs typically fail in this way. See Figure S1 for example of a failed run for YRI. For the remaining runs, we see a range of final classifications accuracies. For the five-parameter models, in YRI this range was 0.5–0.77 (mean 0.619) and for CHB this range was 0.49–0.67 (mean 0.564). To pick the final result, we use the accuracy on the generated data only (i.e. not including the training data); 0.64–1.0 with mean 0.742 for YRI and 0.54–0.9 with mean 0.707 for CHB.

Next, we ran *pg-gan* on 1000 Genomes data from two populations. To model the split of African and non-African populations, we use two pairs of populations separately: YRI/CEU and YRI/CHB, using the OOA2 model from Figure 3c. We use CEU/CHB with the POST model from Figure 3d to represent the post-out-of-Africa split between the ancestors of Europeans and East Asians. The resulting classification accuracies are shown in Figure 8b. The YRI/CEU and YRI/CHB results are comparable to the single-population analysis, but the CEU/CHB classification accuracy is much higher. For all pairs of populations, we provide the parameter inference results in Table 4. Summary statistics for the YRI/CEU split (Figure 9) match the real data closely. YRI/CHB statistics are shown in Figure S4—for the YRI samples these statistics are not quite as closely matched, consistent with the slightly higher classification accuracy for this scenario. CEU/CHB statistics are shown in Figure S5 and are less well-matched to the real data, consistent with the relatively high classification accuracy and suggesting that this model does not contain all the important features for these population, for example archaic admixture or exponential growth.

We also ran *fastsimcoal* on the joint SFS from YRI/CEU, YRI/CHB (using the OOA2 model from Figure 3c), and CEU/CHB (using the POST model from Figure 3d). We used the inferred parameters to create new simulations (with the same fixed mutation rate and HapMap recombination rate distribution used for *pg-gan*). The resulting summary statistics for YRI/CEU are shown in Figure 10, demonstrating that *fastsimcoal* also matches the real data very well. The other *fastsimcoal* results are shown in Figure S6 (YRI/CHB) and Figure S7 (CEU/CHB). For YRI/CHB, *fastsimcoal* produces a slightly better fit than *pg-gan*, but for CEU/CHB the two methods produce very different parameter estimates and neither method matches the summary statistics very well, supporting the suggestion that the generative model is missing some key features of the data.

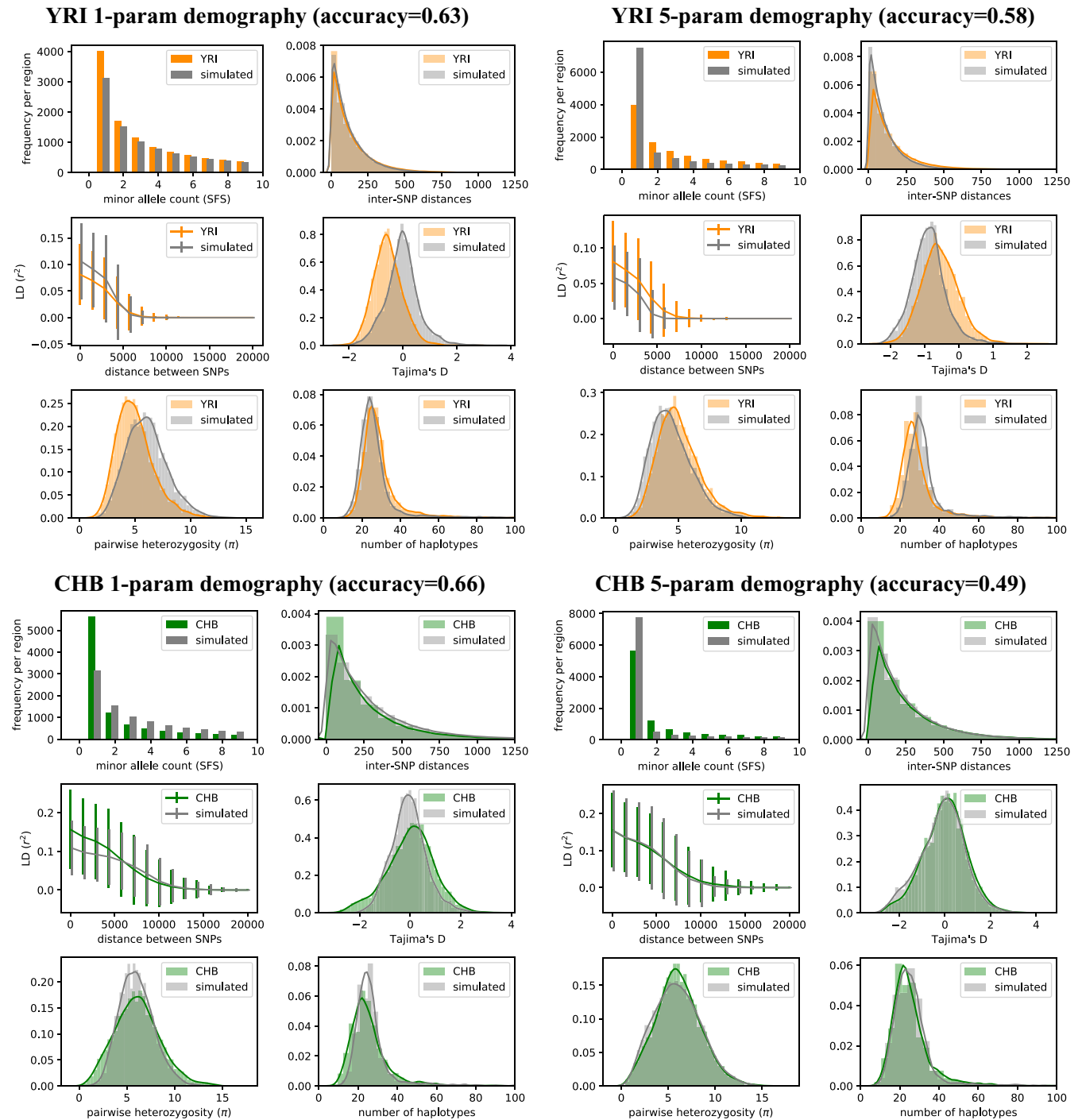


FIGURE 7 Single-population model. Summary statistic comparisons between 1000 Genomes Project data and data simulated under our *pg-gan* inferred parameters for a variety of scenarios. Top left: YRI vs. data simulated under the one-parameter constant population size model. Simulated accuracy: 0.52, overall accuracy: 0.63. Top right: YRI vs. data simulated under the five-parameter exponential growth model. Simulated accuracy: 0.72, overall accuracy: 0.58. Bottom left: CHB vs. data simulated under the one-parameter constant population size model. Simulated accuracy: 0.68, overall accuracy: 0.66. Bottom right: CHB vs. data simulated under the five-parameter exponential growth model. Simulated accuracy: 0.54, overall accuracy: 0.49

Finally, we ran all three populations through the OOA3 model, which was originally described in (Gutenkunst et al., 2009) and recently implemented in *stdpopsim* (Adrion et al., 2020). This required using a 3-population CNN discriminator, which contains many more weights to optimize relative to the two population CNN. In addition,

the OOA3 model requires 14 parameters. We inferred 10 of these parameters, fixing the four migration rate parameters and running *pg-gan* for 500 iterations. We also changed the mutation rate from 2.35×10^{-8} (which was used in (Gutenkunst et al., 2009)) to 1.29×10^{-8} (the recommended human mutation rate from (Adrion

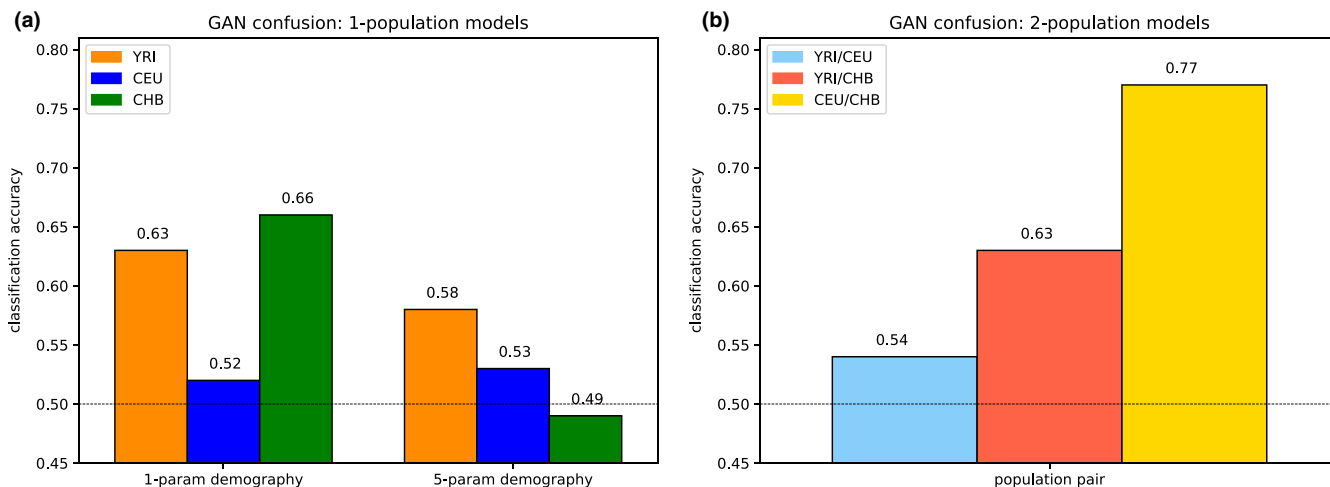


FIGURE 8 GAN confusion for 1- and 2-population models. (a) Comparison of one- and five-parameter models. We use a constant population size for the first group of bars, then move to the five-parameter exponential growth model (Figure 3b). We sample recombination rates from HapMap in both scenarios, instead of fixing the recombination rate. (b) Classification accuracy results on the population split models for YRI/CEU, YRI/CHB and CEU/CHB. The Out-of-Africa models and parameter inference for YRI/CEU and YRI/CHB generally seem to do well, but the CEU/CHB split model and/or parameter inference does not result in simulated data that matches real data

TABLE 3 1000 Genomes single-population parameter inference

Population	N_1	N_2	growth	T_1	T_2
YRI	23,676	20,837	0.0379	2498	1214
CEU	25,127	4676	0.0061	1673	949
CHB	21,136	3150	0.0242	2584	645

Inferred parameters for the exponential growth model (see Figure 3b) in YRI, CEU and CHB. We generally infer similar parameters for CEU and CHB.

et al., 2020)). The inferred parameters and summary statistics are shown in Figure S8, as well as a diagram of the demographic model (reproduced from Figure 2b in (Gutenkunst et al., 2009)). We find a discriminator accuracy of 0.65 and fit some but not all statistics well, suggesting model misspecification, or the difficulty of exploring a relatively high-dimensional parameter space.

3.3 | Computational resources

The runtime of our method is around 5–6 h using a Quadro P5000 GPU. Pre-processing the real data takes several hours for each set of populations (YRI, CEU, CHB, YRI/CEU, YRI/CHB, CEU/CHB and YRI/CEU/CHB). The resulting file sizes are 540 M–944 M, but these do not need to be loaded into memory due to the HDF5 format. The runtime for fastsimcoal was around 55 min.

4 | DISCUSSION

We present a method for automatically learning parameters that can be used to simulate realistic genetic data. Most existing methods

optimize parameters to match summary statistics like the SFS. Our algorithm, *pg-gan*, is a more holistic approach, which finds parameters that generate data that are systematically indistinguishable from the input data, although in practice it also often matches the summary statistics.

Our generative adversarial framework simultaneously trains a generator to produce reasonable evolutionary parameters and a discriminator to distinguish real data from simulated. We use real data during training to make sure the simulations capture realistic genomic features. We demonstrate the use of our method in an isolation-with-migration simulation setting and create simulated data that mirrors three human populations individually, in pairs, and all together. The discriminator often achieves accuracy between 50% and 70%, indicating strong, albeit incomplete, confusion between the real and simulated data. The approach is highly flexible and can automatically fit any parameterized model to any genomic data. We anticipate it will be particularly useful for understudied populations or species, since any unknown parameters can be included in the model and learned.

Our approach yields a natural way of evaluating and refining simulation pipelines. If simulations are easily distinguished from real data, then the model is not producing realistic data. We easily reach essentially complete (50%) discriminator confusion and good summary statistic matching in simulations. But with real data, the fit is imperfect. This could be because there are features of the real data that our models do not include, for example false negatives and other genotyping errors, phasing errors, missing data and inaccessible regions of the genome. Through changes to the generative model, it would be possible to incorporate these effects and evaluate their impact. To handle limited power to detect rare variants (likely why we see more singletons and rare variants in the simulations than the real data), we experimented with filtering a fraction of singletons

TABLE 4 1000 Genomes two-population parameter inference

Populations	N_{anc}	mig	N_1	N_2	N_3	T_1	T_2
YRI/CEU	18,693	-0.0627	4030	27,213	29,863	3501	1132
fastsimcoal	21,017	0.0342	3106	21,954	33,078	2844	1042
YRI/CHB	23,916	0.0738	2422	25,228	27,375	3036	529
fastsimcoal	20,950	0.0167	2959	31,871	32,511	2948	863
CEU/CHB	19,688	-0.0350	16,313	6613	9092	4733	966
fastsimcoal	17,761	0.0240	4044	11,405	15,675	3695	2336

Inferred parameters for the OOA2 model (see Figure 3c) fit to YRI/CEU and YRI/CHB, as well as the POST model (see Figure 3d) fit to CEU/CHB. Results for both *pg-gan* and *fastsimcoal* are included. For *pg-gan*, we generally see similar results for YRI/CEU and YRI/CHB, with a lower classification accuracy for YRI/CEU, indicating a closer match to the real data. Our results are broadly consistent with *fastsimcoal*, except for CEU/CHB, where neither method produces statistics that match the real data (see Figures S5 and S7).

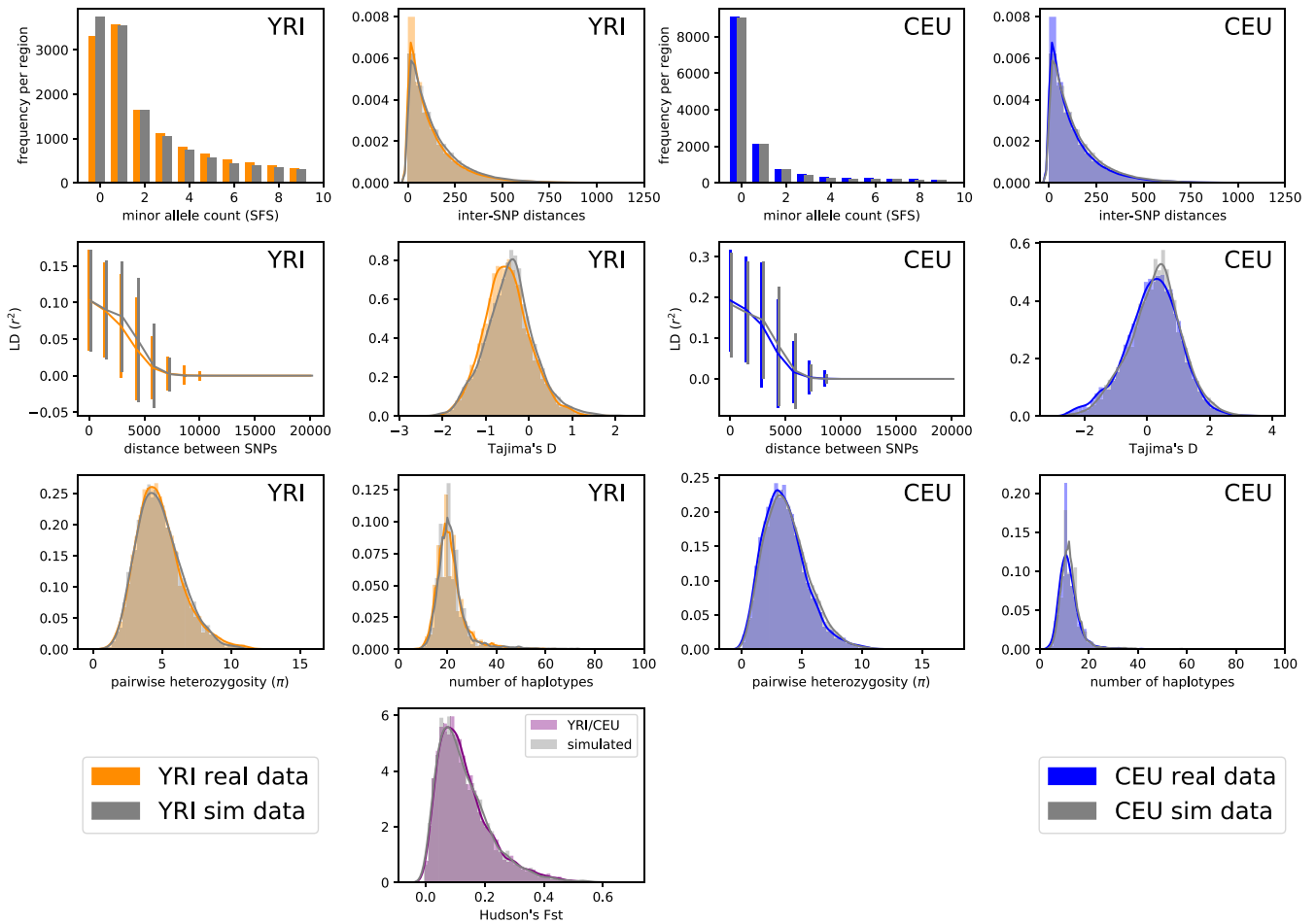


FIGURE 9 YRI/CEU: two-population model. Summary statistic comparison real 1000 Genomes data and data simulated under the inferred parameters from Table 4 (first row). Left: statistics computed on YRI samples only. Right: statistics computed on CEU samples only. Sites with count zero are segregating in only one population. F_{st} between the two populations is shown in the bottom panel. Simulated accuracy: 0.68, overall accuracy: 0.54

from the simulations. This improved the results for YRI, but not for CEU or CHB. Such filtering could be more adaptive in a future iteration. Features such as missing data could be important in some contexts (see *ReLERNN* (Adrien et al., 2020) for example of how to handle missing data). In general, such data quality-related features

are dangerous for our approach, because they provide a way for the discriminator to easily distinguish real and simulated data. For example, if the generative model had data missing at random but the real data are missing in a non-random fashion, then the discriminator will use this signal for classification. It would be important to make the

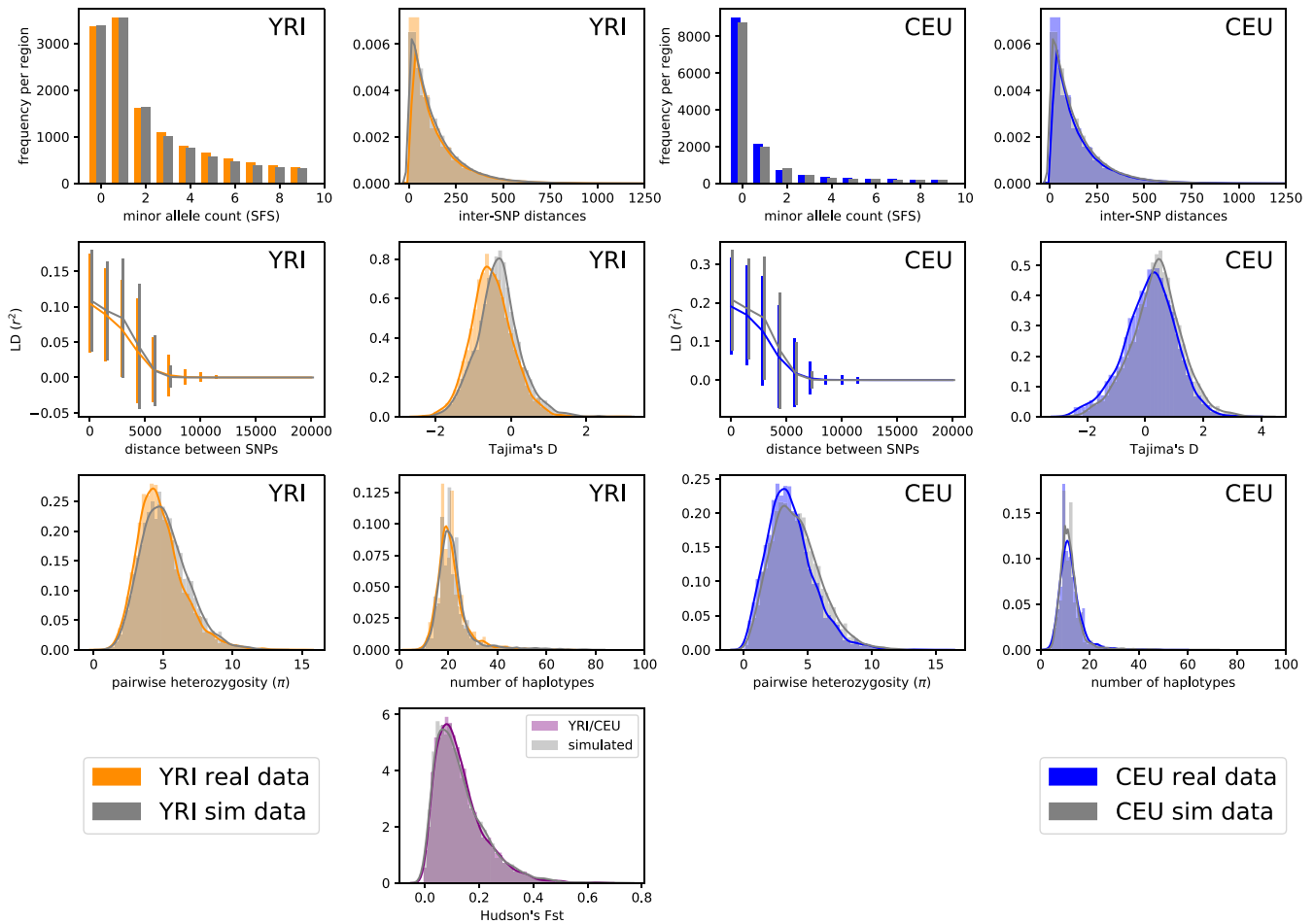


FIGURE 10 YRI/CEU: two-population model (*fastsimcoal*). Summary statistic comparison between YRI/CEU and data simulated under the OOA2 model parameters inferred by *fastsimcoal*. Here, we include all the statistics (unlike Figure 6) since we are providing *fastsimcoal* with a recombination rate distribution. Left: statistics computed on YRI samples only. Right: statistics computed on CEU samples only. Sites with count zero are segregating in only one population. F_{st} between the two populations is shown in the bottom panel

generative model sufficiently flexible that it could learn to replicate the distinguishing features of the real data.

Some subsets of populations were more difficult to fit than others. The CEU/CHB split proved particularly difficult for both *pg-gan* and *fastsimcoal*. Since data quality should be similar between populations, this probably indicates that our model does not include demographic features that are important for patterns of variation in these populations, for example archaic admixture or exponential growth. More generally, our model ignores many important biological features, for example heterogeneity in the mutation rate and other parameters, and natural selection. We assumed that mutation and recombination rates were known, but they can easily be added as parameters to the generative model and inferred (Figure S9). Heterogeneity could be modelled by fitting a distribution from which to draw parameters, rather than a point estimate. Natural selection, which can bias estimates of demographic parameters (Schridder et al., 2016), is more difficult to model. The effects of regions under strong positive selection or long-term balancing selection can be minimized by removing them from the training data. However, background

selection affects the majority of the genome and completely restricting to 'truly neutral' regions of the genome is impractical. One simple but somewhat unsatisfactory solution is to approximate the effects of background selection across the genome by scaling effective population sizes with a factor drawn from empirical estimates of the effect of background selection across the genome (McVicker et al., 2009). A better solution would be to estimate the distribution of selection coefficients as part of the model (Johri et al., 2021). This requires a generator that can simulate selection, for example *SLiM* (Haller & Messer, 2019) but would be much more computationally intensive than the coalescent simulations in the current approach. Efficiently incorporating selection into the model is a key area for future development.

There are several areas of future exploration that involve algorithmic modifications. In our current implementation, the topology of the demographic model needs to be specified ahead of time. However, it would be possible to extend our method to explore a space of demographic models, which would allow both the topology and the model parameters to be learned automatically. Although

we mitigate overfitting by selecting real data regions at random (as opposed to a fixed sliding window), it is still a concern for the discriminator due to the fundamental data imbalance. The amount of real data is fixed, but the number of simulated examples is unlimited. There are many ways to guard against overfitting neural networks, including regularization and architecture modifications. An important line of future research is to optimize the training procedure in the presence of limited real data.

Another asymmetry comes from the potentially different learning rates of the generator and discriminator. The training of both components needs to be balanced—if the discriminator learns the difference between real and simulated data too quickly, the generator might not have a chance to explore a parameter space that would actually cause confusion. On the other hand, if the discriminator learns too slowly, all generator updates might look equally confusing. It would be interesting to explore adaptively controlling the learning rate—slowing down either the generator or the discriminator as needed through fewer parameter proposals or mini-batches. Understanding the behaviour of the discriminator is itself an important area of future work, which could help us investigate alignment between its hidden layers and traditional summary statistics.

Some idea of the uncertainty in the parameter space can be obtained by looking at the distribution of replicate estimates. In principle, this approach could be extended to provide bootstrap confidence intervals by fitting the model to resampled data. A more general approach would be to fix the discriminator and vary the generator parameters to identify the parameter space over which the discriminator has low accuracy.

Our approach could be incorporated into a transfer learning (Pan & Yang, 2010) framework. In transfer learning, the parameters of an ML model are initialized by training on a large data set, then 'fine-tuned' by training on a smaller number of examples from the target data set. In our case, a large data set like the 1000 Genomes could be used to find an initial guess for discriminator weights, then these weights could be fine-tuned using data with fewer regions or sequenced individuals. The evolutionary model could still be modified, as transfer learning would be used for the discriminator, not the generator. The only restriction would be that the number of populations would need to match between the larger data set and the smaller data set. The original learning on the larger data set would primarily assist the discriminator in learning general features of genomic data sets—population-level specifics would be learned in the fine-tuning phase.

It is our hope that others will build upon this initial exploration into parametric GANs for population genetics. Future developments will include integrating more realistic features of real data, constructing bootstrap confidence intervals for parameter estimates, and applying our approach to non-human species. In terms of methodological development, we aim to integrate transfer learning and develop interpretative approaches for the CNN discriminator, in order to investigate alignment between its hidden layers and traditional summary statistics. Modern machine learning has proved to be powerful in many domains, and our work emphasizes that this

is true for population genetics as well. However, machine learning in population genetics requires novel architectures, for example our parametric generator and multi-population CNN discriminator—innovations that will be useful for future development of ML methods in the field.

ACKNOWLEDGEMENTS

The authors would like to thank Joe Cammisa for extensive computational support and Ke Wang for assistance with comparing *pg-gan* to other methods. SM is funded in part by NIH grant R15HG011528, and IM is funded in part by NIH grant R35GM133708. MK was funded through a David Robbins '83 Big Data/Social Change Lang Center Internship. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health or other funding sources.

DATA AVAILABILITY STATEMENT

Our *pg-gan* software uses a *tensorflow* (Abadi et al., 2015) backend and is available open-source at <https://github.com/mathiesonlab/pg-gan>. All data included in this work are publicly available through the 1000 Genomes Project <https://www.internationalgenome.org/> (1000 Genomes Project Consortium, 2015).

ORCID

Sara Mathieson  <https://orcid.org/0000-0002-0484-0838>

REFERENCES

- 1000 Genomes Project Consortium. (2015). A global reference for human genetic variation. *Nature*, 526(7571), 68–74.
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., & Jia, Y., ... Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Available from: <https://www.tensorflow.org/Software/tensorflow.org>
- Adrion, J. R., Cole, C. B., Dukler, N., Galloway, J. G., Gladstein, A. L., Gower, G., Kyriazis, C. C., Ragsdale, A. P., Tsambos, G., Baumdicker, F., Carlson, J., Cartwright, R. A., Durvasula, A., Gronau, I., Kim, B. Y., McKenzie, P., Messer, P. W., Noskova, E., Ortega-Del Vecchyo, D., ... Kern, A. D. (2020). A community-maintained standard library of population genetic models. *eLife*, 9, <https://doi.org/10.7554/eLife.54967>
- Adrion, J. R., Galloway, J. G., & Kern, A. D. (2020). Predicting the landscape of recombination using deep learning. *Molecular Biology and Evolution*, 37(6), 1790–1808.
- Batthey, C. J., Coffing, G. C., & Kern, A. D. (2021). Visualizing population structure with variational autoencoders. *G3 Genes, Genomes, Genetics*, 11(1), 1–11. <https://doi.org/10.1093/g3journal/jkaa036>
- Beaumont, M. A., Zhang, W., & Balding, D. J. (2002). Approximate Bayesian computation in population genetics. *Genetics*, 162(4), 2025–2035.
- Beichman, A. C., Phung, T. N., & Lohmueller, K. E. (2017). Comparison of single genome and allele frequency data reveals discordant demographic histories. *G3: Genes, Genomes, Genetics*, 7(11), 3605–3620.
- Blum, M. G. B., & François, O. (2010). Non-linear regression models for Approximate Bayesian Computation. *Statistics and Computing*, 20(1), 63–73.

- Borji, A. (2019). Pros and cons of gan evaluation measures. *Computer Vision and Image Understanding*, 179, 41–65.
- Chan, J., Perrone, V., Spence, J., Jenkins, P., Mathieson, S., & Song, Y. (2018). A likelihood-free inference framework for population genetic data using exchangeable neural networks. *Advances in Neural Information Processing Systems*, 8594–8605.
- Dieng, A. B., Ruiz, F. J. R., Blei, D. M., & Titsias, M. K. (2019). Prescribed generative adversarial networks. *arXiv preprint arXiv:1910.04302*.
- Ewing, G., & Hermisson, J. (2010). MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics*, 26(16), 2064–2065.
- Excoffier, L., Dupanloup, I., Huerta-Sánchez, E., Sousa, V. C., & Foll, M. (2013). Robust demographic inference from genomic and SNP data. *PLoS Genetics*, 9(10), e1003905.
- Excoffier, L., & Foll, M. (2011). Fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics*, 27(9), 1332–1334.
- Flagel, L., Brandvain, Y., & Schrider, D. R. (2019). The unreasonable effectiveness of convolutional neural networks in population genetic inference. *Molecular Biology and Evolution*, 36(2), 220–238.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2672–2680.
- Gower, G. R., Picazo, P. I., Fumagalli, M., & Racimo, F. (2020). Detecting adaptive introgression in human evolution using convolutional neural networks. *bioRxiv*.
- Griffiths, R. C., & Marjoram, P. (1997). An ancestral recombination graph. *IMA*, 87, 257.
- Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H., & Bustamante, C. D. (2009). Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genetics*, 5(10), e1000695.
- Haller, B. C., & Messer, P. W. (2019). SLiM 3: Forward genetic simulations beyond the Wright-Fisher model. *Molecular Biology and Evolution*, 36(3), 632–637.
- Harris, K. (2015). Evidence for recent, population-specific evolution of the human mutation rate. *Proceedings of the National Academy of Sciences*, 112(11), 3439–3444.
- Harris, K., & Pritchard, J. K. (2017). Rapid evolution of the human mutation spectrum. *Elife*, 6, e24284.
- Hinch, A. G., Tandon, A., Patterson, N., Song, Y., Rohland, N., Palmer, C. D., Chen, G. K., Wang, K., Buxbaum, S. G., Akylbekova, E. L., Aldrich, M. C., Ambrosone, C. B., Amos, C., Bandera, E. V., Berndt, S. I., Bernstein, L., Blot, W. J., Bock, C. H., Boerwinkle, E., ... Myers, S. R. (2011). The landscape of recombination in African Americans. *Nature*, 476(7359), 170–175. <https://doi.org/10.1038/nature10336>
- Hudson, R. R. (2002). Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*, 18(2), 337–338.
- Hudson, R. R., Slatkin, M., & Maddison, W. P. (1992). Estimation of levels of gene flow from DNA sequence data. *Genetics*, 132(2), 583–589.
- International HapMap Consortium, et al. (2007). A second generation human haplotype map of over 3.1 million SNPs. *Nature*, 449(7164), 851.
- Johri, P., Riall, K., Becher, H., Excoffier, L., Charlesworth, B., & Jensen, J. D. (2021). The impact of purifying and background selection on the inference of population history: problems and prospects. *bioRxiv*, <https://doi.org/10.1101/2020.04.28.066365>. URL: <https://www.biorxiv.org/content/early/2021/01/18/2020.04.28.066365>
- Kelleher, J., Etheridge, A. M., & McVean, G. (2016). Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Computational Biology*, 12(5), e1004842.
- Kern, A. D., & Schrider, D. R. (2016). Discoal: flexible coalescent simulations with selection. *Bioinformatics*, 32(24), 3839–3841.
- Kessler, M. D., Loesch, D. P., Perry, J. A., Heard-Costa, N. L., Taliun, D., Cade, B. E., Wang, H., Daya, M., Ziniti, J., Datta, S., Celedón, J. C., Soto-Quiros, M. E., Avila, L., Weiss, S. T., Barnes, K., Redline, S. S., Vasani, R. S., Johnson, A. D., Mathias, R. A., ... O'Connor, T. D. (2020). De novo mutations across 1,465 diverse genomes reveal mutational insights and reductions in the Amish founder population. *Proceedings of the National Academy of Sciences*, 117(5), 2560–2569. <https://doi.org/10.1073/pnas.1902766117>
- Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., & Yosef, N. (2018). Deep generative modeling for single-cell transcriptomics. *Nature Methods*, 15(12), 1053–1058.
- Mathieson, I., & Reich, D. (2017). Differences in the rare variant spectrum among human populations. *PLoS Genetics*, 13(2), e1006581.
- McVicker, G., Gordon, D., Davis, C., & Green, P. (2009). Widespread genomic signatures of natural selection in hominid evolution. *PLoS Genetics*, 5(5), e1000471. <https://doi.org/10.1371/journal.pgen.1000471>
- Miles, A. (2015). Estimating Fst. <http://alimanfoo.github.io/2015/09/21/estimating-fst.html>
- Miles, A. (2017). Extracting data from VCF files. <http://alimanfoo.github.io/2017/06/14/read-vcf.html>
- Mirza, M., & Osindero, S. (2014). Conditional generative adversarial nets. *arXiv Preprint arXiv:1411.1784*.
- Neuhauser, C., & Krone, S. M. (1997). Ancestral processes with selection. *Theoretical Population Biology*, 51, 210–237.
- Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345–1359.
- Pincus, M. (1970). Letter to the editor—a Monte Carlo method for the approximate solution of certain types of constrained optimization problems. *Operations Research*, 18(6), 1225–1228.
- Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Riesselman, A. J., Ingraham, J. B., & Marks, D. S. (2018). Deep generative models of genetic variation capture the effects of mutations. *Nature Methods*, 15(10), 816–822.
- Ronen, R., Udpa, N., Halperin, E., & Bafna, V. (2013). Learning natural selection from the site frequency spectrum. *Genetics*, 195(1), 181–193.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., & Chen, X. (2016). Improved techniques for training gans. *arXiv preprint arXiv:1606.03498*
- Sanchez, T., Cury, J., & Charpiat, G., & Jay, F. (2020). Deep learning for population size history inference: Design, comparison and combination with approximate Bayesian computation. *Molecular Ecology Resources*. <https://doi.org/10.1111/1755-0998.13224>
- Schrider, D. R., Shanku, A. G., & Kern, A. D. (2016). Effects of linked selective sweeps on demographic inference and model selection. *Genetics*, 204(3), 1207–1223. <https://doi.org/10.1534/genetics.116.190223>
- Sheehan, S., & Song, Y. S. (2016). Deep learning for population genetic inference. *PLOS Computational Biology*, 12(3), e1004845. <https://doi.org/10.1371/journal.pcbi.1004845>
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1), 1929–1958.
- Terhorst, J., Kamm, J. A., & Song, Y. S. (2017). Robust and scalable inference of population history from hundreds of unphased whole genomes. *Nature Genetics*, 49(2), 303–309.
- Teshima, K. M., & Innan, H. (2009). mbs: modifying Hudson's ms software to generate samples of DNA sequences with a biallelic site under selection. *BMC Bioinformatics*, 10(1), 166.
- Torada, L., Lorenzon, L., Beddis, A., Isildak, U., Pattini, L., Mathieson, S., & Fumagalli, M. (2019). ImaGene: A convolutional neural network to quantify natural selection from genomic data. *BMC Bioinformatics*, 20(9), 337.

- Xu, Q., Huang, G., Yuan, Y., Guo, C., Sun, Y., Wu, F., & Weinberger, K. (2018). An empirical study on evaluation metrics of generative adversarial networks. *arXiv preprint arXiv:1806.07755*.
- Yelmen, B., Decelle, A., Ongaro, L., Marnetto, D., Tallec, C., Montinaro, F., Furtlehner, C., Pagani, L., & Jay, F. (2021). Creating artificial human genomes using generative models. *PLoS Genetics*, 17(2), e1009303.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

How to cite this article: Wang Z, Wang J, Kourakos M, et al. Automatic inference of demographic parameters using generative adversarial networks. *Mol Ecol Resour*. 2021;21:2689–2705. <https://doi.org/10.1111/1755-0998.13386>