

Simultaneous Inference of Past Demography and Selection from the Ancestral Recombination Graph under the Beta Coalescent

Kevin Korfmann^{1*}, Thibaut Sellinger^{1*,2*}, Fabian Freund^{3,5},
Matteo Fumagalli⁴, Aurélien Tellier¹⁺

¹ Population Genetics, Department of Life Science Systems,
Technical University of Munich (Liesel-Beckmann-Strasse 2, 85354 Freising, Germany)

² Department of Environment and Biodiversity, Paris Lodron University of Salzburg

³ Institute of Plant Breeding, Seed Science and Population Genetics, University of Hohenheim
(Fruwirthstrasse 21, 70599 Stuttgart, Germany)

⁴ Department of Biological and Behavioural Sciences, Queen Mary University of London
(Mile End Road, London E1 4NS, UK)

⁵ Department of Genetics and Genome Biology, University of Leicester
(University Road, Leicester LE1 7RH, UK)

* both author contributed equally and share first authorship

+ Corresponding author, aurelien.tellier@tum.de

Abstract

The reproductive mechanism of a species is a key driver of genome evolution. The standard Wright-Fisher model for the reproduction of individuals in a population assumes that each individual produces a number of offspring negligible compared to the total population size. Yet many species of plants, invertebrates, prokaryotes or fish exhibit neutrally skewed offspring distribution or strong selection events yielding few individuals to produce a number of offspring of up to the same magnitude as the population size. As a result, the genealogy of a sample is characterized by multiple individuals (more than two) coalescing simultaneously to the same common ancestor. The current methods developed to detect such multiple merger events do not account for complex demographic scenarios or recombination, and require large sample sizes. We tackle these limitations by developing two novel and different approaches to infer multiple merger events from sequence data or the ancestral recombination graph (ARG): a sequentially Markovian coalescent (SM β C) and a graph neural network (GNN $_{coal}$). We first give proof of the accuracy of our methods to estimate the multiple merger parameter and past demographic history using simulated data under the β -coalescent model. Secondly, we show that our approaches can also recover the effect of positive selective sweeps along the genome. Finally, we are able to distinguish skewed offspring distribution from selection while simultaneously inferring the past variation of population size. Our findings stress the aptitude of neural networks to leverage information from the ARG for inference but also the urgent need for more accurate ARG inference approaches.

Keywords— kingman coalescent, beta coalescent, selective sweep, deep learning, graph neural networks, population genetics, multiple merger coalescent, sequentially markovian coalescent, ancestral recombination graph

Introduction

With the availability of genomes of increasing quality for many species across the tree of life, population genetics models and statistical methods have been developed to recover the past history of a population/species from whole genome sequence data from several individuals [83, 54, 78, 84, 81, 5, 4, 86, 41, 42]. Indeed, the inference of the past demographic history of a species, *i.e.* population expansion, contraction, or bottlenecks, extinction/colonisation, is not only interesting in its own right, but also essential to calibrate genome-wide scans to detect genes under (*e.g.* positive or balancing) selection [86, 43]. A common feature of inference methods that make full use of whole genome sequences is the underlying assumption of a Kingman coalescent process [48] to describe the genealogy distribution of a sample. The Kingman coalescent process and its properties stem from using the traditional forward-in-time Wright-Fisher (WF) model to describe the reproduction mechanism of a population. Besides non-overlapping generations, a key assumption of the neutral WF model is that an individual offspring chooses randomly (*i.e.* uniformly) its parents from the previous generation. More precisely, each chromosome chooses a parental chromosome from the previous generation. Thus, a key parameter is the distribution of the number of offspring that parents can have. In the WF model, due to the binomial sampling, the distribution of offspring number per parent is well approximated by a Poisson distribution with both mean and variance equal to one. This implies that parents will most likely have zero, one, or two offspring individuals, but it is improbable that one parent would have many offspring individuals (*i.e.* on the order of the population size, under the Wright-Fisher haploid model the probability for a parent to have 10 or more offspring is $\approx 10^{-8}$). The assumption of small variance in offspring distribution between individual parents is realistic for species with low juvenile mortality (so-called type I and II survivorship in ecology), such as mammals.

As genome sequence data become available for a wide variety of species with different biological traits and/or life cycles, the applicability of the Kingman coalescent relying on the WF model can be questioned [85, 2, 3, 65, 44, 62, 88, 59, 30]. Indeed, for some species, such as fish, with high fecundity and high juveniles mortality (type III survivorship), it is expected that the variance in reproduction between parents can be much larger than under the Poisson distribution [88]. This effect is termed as sweepstake reproduction [35, 2]. Neutral processes such as strong seed banking [12], high fecundity with skewed offspring distribution [35, 25], extremely strong and recurrent bottlenecks [9, 21], and strong selective processes (*i.e.* positive selection) [24, 17, 18, 34, 3] are theoretically shown to deviate from the classic WF model in a way that the genealogies can no longer be described by a Kingman coalescent process. Under such conditions, a new class of processes arise to describe the genealogy distribution, a class where multiple individuals can coalesce and/or multiple distinguished coalescence events can occur simultaneously [74, 61, 23, 73, 67, 14]. Generally, this class of genealogical processes is called the Multiple Merger Coalescent (MMC). MMC models are more biologically appropriate than the Kingman coalescent to study many species of fish [26, 2, 3, 35], invertebrates (insects, crustaceans, etc.), viruses [57], bacteria [59, 63], plants and their pathogens [88]. While we would like to assess which population model best describes the species genealogy, field experiments to quantify the underlying reproduction mechanism of a species can be costly and time consuming at best, or intractable at worst. Therefore, an alternative solution is to use inference methods based on genome data to identify which model best describes

Inference under the Beta Coalescent

the genealogy of a given species/population.

In this study we use the so-called β -coalescent, a specific class of MMC models. Unlike under the WF model, under MMC models the ploidy level strongly affects the distribution of genealogies [8]. For simplicity, in this study we focus on haploid organisms. It is demonstrated that if the probability of a parent to have k or more offspring is proportional to $k^{-\alpha}$, where $1 < \alpha < 2$, then the genealogy can be described by a Λ -coalescent [80]. The latter is a general class of coalescent process describing how and how fast ancestral lineages merge [67, 73]. When using the Beta($2-\alpha, \alpha$) distribution as a probability measure for the Λ -coalescent, the transition rates (*i.e.* coalescent rate) can be analytically obtained leading to the β -coalescent, a specific MMC model. If α tends to 2, then the coalescent process converges to a Kingman coalescent (up to a scaling constant) [8, 51, 52]. If α tends to one, the model tends to a Bolthausen-Sznitman coalescent process (*i.e.* dominated by strong multiple merger events) [14]. The β -coalescent has the property that the observed polarized Site Frequency Spectrum (SFS) of a sample of single nucleotide polymorphisms (SNPs) exhibits a characteristic U-shape with an excess of rare and high frequency variants (compared to the Kingman coalescent) [77]. Current methods to draw inference under MMC models leverage information from the summary statistics extracted from full genome data such as Site Frequency Spectrum (SFS, or derived summary statistics) [52, 34, 72], minor allele frequency [70] or copy number alteration [44]. It is shown that the SFS is robust to the effect of recombination [52, 70] and its shape allows to discriminate between simple demographic models (population expansion or contraction) under the Kingman coalescent and MMC models with constant population size [52, 51, 26]. However, methods relying on genome-wide SFS have two main disadvantages. First, in absence of strong prior knowledge, they can suffer from non-identifiability [41] as several complex neutral demographic and/or selective models under the Kingman or MMC models can generate similar SFS distributions. Second, as they summarize the collection of underlying genealogies, they require high sample sizes (>50) to produce trustworthy results [52, 51, 26], relying on experimental designs which are prohibitive for the study of non-model species. To tackle these limitations, we develop two methods that integrate recombination events along the genome in order to leverage more information from full genome data, thus requiring fewer samples.

In species undergoing sexual reproduction, recombination events break the genealogy of a sample at different position of the genome (*i.e.* the genealogy of a sample varies along the genome), leading to what is called the Ancestral Recombination Graph (ARG) [38, 8]. Because all the genealogical information is contained in the ARG, in this study we aim at the interpretation of the ARGs to recover model parameters in presence of multiple merger events. With the development of the sequentially Markovian coalescent theory [58, 56, 93], it becomes tractable to integrate linkage disequilibrium over chromosomes in inferences based on the Kingman coalescent [54]. Hence, we first develop an SMC approach based on the β -coalescent named the Sequentially Markovian β Coalescent (SM β C). The β -coalescent has the additional property that, under recombination, long range dependency can be generated between coalescent trees along the genome if multiple-merger events happen in a single generation [8]. In other words, coalescent trees which are located at different places in the genome, and expected to be unlinked from one another [64], would show non-zero correlation in their topology and coalescent times. This is because coalescent trees from different genomic regions may all be affected by the

same MMC event (merger event of multiple lineages in the past) which then leaves traces in the genome at several loci [9]. To overcome the theoretically predicted non-Markovian property of the distribution of genealogies along the genome under the β -coalescent with recombination [8], we develop a second method based on deep learning (DL) trained from efficient coalescent simulations [7]. In evolutionary genomics, DL approaches trained by simulations are shown to be powerful inference tools [83, 50]. Previous work demonstrated that DL approach can help overcome problems mathematically insolvable or computationally intractable in the field of population genetics [83, 6, 92, 96, 29, 22, 68, 19, 40]. The novelty of our neural network relies on its structure (Graph Neural Network, GNN) and its training algorithm based on the ARG of a sample, or its tree sequence representation [45]. GNNs are an emerging category of DL algorithm [16, 94, 20, 99] that benefit by using irregular domain data (*i.e.* graphs). GNNs are designed for the prediction of node features [49, 95], edge features (link prediction) [98, 79], or additional properties of entire graphs [97, 53]. Therefore, GNNs represent a new tool to address the large dimensionality of ARGs, while simultaneously leveraging information from the genealogy (namely topology and age of coalescent events) as a substantial improvement over convolutions of genotype matrices, as currently done in the field [75].

We first quantify the bias of previous SMC methods (MSMC and MSMC2 [78, 91]) when performing inference of past population size variation under the β -coalescent. We then describe our two methods, $SM\beta C$ and GNN_{coal} , and demonstrate their statistical power as well as their respective limitations. From simulated tree-sequence (*i.e.* ARG) and sequence (*i.e.* SNPs) data, we assess the accuracy of both approaches to recover the past variation of population size and the α parameter of the Beta-distribution. This parameter indicates how frequent and strong multiple merger events occur. We demonstrate that our approaches can infer the evolutionary mechanism responsible for multiple merger events and distinguish local selection events from genome-wide effects of multiple mergers. We highlight the limits of the Markovian property of SMC to describe data generated under the β -coalescent. Finally, we show that both our approaches can model and identify the presence of selection along the genome while simultaneously accounting for non-constant population size, recombination, and skewed offspring distribution. Thus our methods represents a major and necessary leap forward in the field of population genetic inferences.

Materials and Methods

In our study we first assume the true ARG to be known. Hence, the ARG of the sample is given as input to our methods to estimate recover model parameters of interest (*e.g.* the α parameter and/or the past variation of population size). We then show the applicability of our methods by using as input simulated sequence data (*i.e.* SNPs) and/or ARG inferred using ARGweaver [69] from simulated sequence data.

SMC-based method

In this study, we use different SMC-based algorithms: two previously published, MSMC and MSMC2 [78, 91], and the new $SM\beta C$. In the latter, the software backbone stems from our previous eSMC [81, 82] whilst the theoretical framework originates from the MSMC algorithm [78] (see Supplementary Text S1). All approaches can use the ARG

Inference under the Beta Coalescent

or sequence data as input. Giving ARG as input for MSMC an MSMC2 is enabled by a re-implementation included in the R package eSMC2 [82]. The MSMC2 algorithm focuses on the coalescence time between two haploid samples along the genome. In the event of recombination, there is a break in the current genealogy and the coalescence time consequently takes a new value. A detailed description of the algorithm can be found in [27, 91]. The MSMC algorithm simultaneously analyses multiple sequences (up to 10) and follows the distribution of the first coalescence event in a sample of size $n > 2$ along the sequence based on the Kingman coalescent [48]. A detailed description of MSMC can be found in [78].

Our new approach, $SM\beta C$, is a theoretical extension of the MSMC algorithm, simultaneously analyzing multiple haploid sequences and focusing on the first coalescence event of a sample size 3 or 4. The $SM\beta C$ follows the distribution of the first coalescence event of a sample along sequences assuming a β -coalescent process. Therefore, our $SM\beta C$ allows for more than two ancestral lineages to join the first coalescence event, or new lineages to join an already existing binary (or triple) coalescent event. Hence, the $SM\beta C$ extends the MSMC theoretical framework by adding hidden states at which more than two lineages coalesce. Currently, the $SM\beta C$ has been derived to analyze for up to 4 sequences simultaneously (due to computational load and mathematical complexity). The emission matrix is similar to the one of MSMC. As in the MSMC software, the population size is assumed piece-wise constant in time and we discretize time in 40 bins throughout this study. A detailed description of $SM\beta C$ can be found in Supplementary Text S1. To test and validate the theoretical accuracy of our approach, we first study its best case convergence (introduced in [82]) which corresponds to the model's performance when the true genealogy is given as input, *i.e.* as if the hidden states are known. Additionally, we also validate the practical accuracy of the $SM\beta C$ on simulated sequence data taking the same input as the MSMC software [78], or using the inferred ARGs by ARGweaver [69]. All SMC approaches used in this manuscript are found in the R package eSMC2 (<https://github.com/TPPSellinger/eSMC2>).

GNN_{coal} method

Inspired by results obtained from inferences based on tree sequence data [32, 82], we develop a graph neural network (GNN) taking tree sequence data as input. Our GNN is designed to infer population size along with the α parameter of the Beta distribution describing the distribution of offspring production. In practice, the ARG is reshaped into a sequence of genealogies (more precisely a sequence of undirected graphs), and then given as input to the GNN (similar to what is described above for the $SM\beta C$). In our analyses, we fixed the batch size to 500. This value represents the number coalescence trees being processed before updating parameters of the neural network. As batch size is fixed to 500, only simulations displaying at least 500 recombination events are considered for the training data sets. If more than 500 recombination events occur along the sequence, the ARG is truncated and the GNN will only take as input the first 500 genealogies and remove the rest. Thanks to the GNN architecture, the algorithm can account for the topology of the genealogy. Hence, the GNN leverages information from coalescence time and branch lengths but also from the topology of the ARG. This operation is known as a graph convolution. By doing so, the GNN is capable of learning from local features of the ARG and extract information from its complex structure. To learn from global ge-

Inference under the Beta Coalescent

nealogy patterns (which SMC-based methods cannot do), an additional pooling strategy is implemented as part of the network [97]. To do so, the ARG is broken into smaller ARGs (*i.e.* subgraphs) during the forward-pass step. To illustrate the GNN strategy, we visualize the compression-like process, from the coalescent trees (1) being processed by GNNcoal (2,3) to the inferred variable of interest (4, 5) in Figure 1.

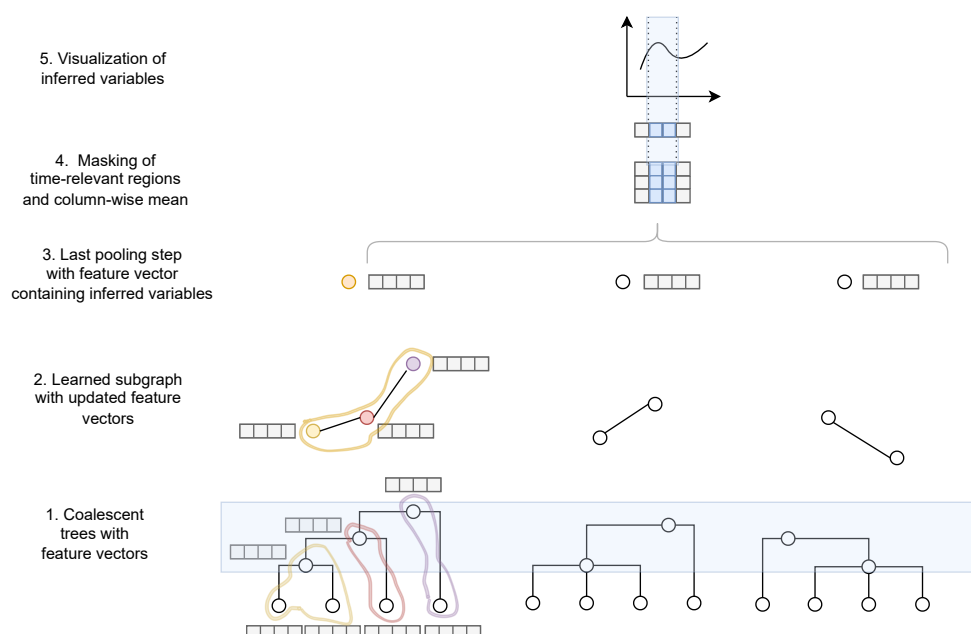


Fig. 1. Schematic representation of GNNcoal processing an ARG Hierarchical pooling compression of a sequence of coalescent trees into the inferred variable of interest (*i.e.* demographic changes) using graph convolutions. Each coalescent ancestor or leaf node is initialized by a feature vector (light grey boxes) (1). Sub-graphs are generated by a pooling network with updated feature vectors and a final compression step is performed until ideally one node per graph remains (2-3). Lastly, the column-wise mean is taken after applying a time mask (blue - based on number of coalescent events), so that single feature vector remains (4-5). Detailed description of the graph convolution, feature vector initialization, pooling methodology, coalescent time mask construction, and dataset generation can be found in Supplementary Text S2.

To infer parameters from our neural network, we need to define an objective function to be optimized. We use a masked root-mean-squared error (RMSE) loss function as objective function which is computed for each inputted ARG (*i.e.* minimizing the average square difference between predicted and true parameter value). In practice, time is discretized (as for the SM β C) and time windows are defined. The true α value and true demography at 60 predefined time points are given as input to the GNN to compute the loss function. The GNN captures the stochastic complexity arising from the underlying demographic scenario and model parameters. Furthermore, our algorithm naturally defines an appropriate time window to have sufficient observation at each time point. A more detailed description of the GNNcoal can be found in Supplementary Text S2. The code of the model architecture is implemented in *Pytorch* [66] using the extension *Pytorch Geometric* [28]. The model is available with the simulated training dataset at <https://github.com/kevinkorfmann/GNNcoal> and <https://github.com/kevinkorfmann/GNNcoal-analysis>.

ARGweaver

As the ARG is not known in practice, it needs to be inferred from sequence data. ARGweaver displays the best performance at recovering the ARG from whole genome polymorphism data at the sample sizes employed in this study (*i.e.* $\ll 50$) [69, 15]. Briefly, ARGweaver samples the ARG of n chromosomes/scaffolds conditional on the ARG of $n - 1$ chromosomes/scaffolds. To this aim, ARGweaver relies on hidden Markov models while assuming a sequentially Markov coalescent process and a discretization of time, similarly to the SMC-based methods previously described. For a more detail description of the algorithm, we refer the reader to the supplementary material of [69].

Simulation of data

Validation dataset for both methods

The ARG is given as input to the DL approach and the SM β C (see [82]). We use msprime [7] to simulate the ARG of a sample (individuals are assumed to be haploid) under the β -coalescent based on [80, 8] or under the Kingman coalescent (under neutrality or selection). We simulate 10 sequences of 100 Mbp under five different demographic scenarios: 1) Constant population size; 2) Bottleneck with sudden decrease of the population size by a factor 10 followed by a sudden increase of population by a factor 10; 3) Expansion with sudden increase of the population size by a factor 10; 4) Contraction with sudden decrease of the population size by a factor 10; and 5) "Saw-tooth" with successive exponential decreases and increases of population size through time, resulting in continuous population size variation (as shown in [89, 78, 82]). We simulate data under different α values (*i.e.* parameters of the β -distribution) including values of 1.9 (almost no multiple merger events), 1.7, 1.5, and 1.3 (frequent and strong multiple merger events). Mutation and recombination rate (respectively μ and r) are set to 10^{-8} per generation per bp in order to obtain the best compromise between realistic values and number of SNPs. When specified, some specific scenarios assume recombination and mutation rate set to produce sufficient data or to avoid violation of the finite site hypothesis. All python scripts used to simulate data sets are available at <https://github.com/kevinkorfmann/GNNcoal-analysis>.

Additionally, to generate sequence data, we simulate 10 sequences of 10 Mbp under the five different demographic scenarios described above and for the same α values. For each scenario, 10 replicates are simulated. In order to obtain sufficient SNPs for inference, we simulate sequence data with mutation and recombination rate (respectively μ and r) of 10^{-8} per generation per bp when α is set to 1.9 and 1.7, 10^{-7} per generation per bp when α is set to 1.5, and 10^{-6} per generation per bp when α is set to 1.3.

Training dataset for the GNN_{coal}

In our study we train two GNNs, one to infer past variation of population size through time along with α , and one for model selection. The training dataset used for both GNNs is described below.

Training dataset for the GNN inferring α and demography

We generate an extensive number of ARGs to train our GNN. The ARGs are simulated

Inference under the Beta Coalescent

under many demographic scenarios and α values. The model parameters are updated in supervised manner. The loss function is calculated for each batch with respect to how much the machine-learning estimates differ from the true parameters used for simulation. The simulations strategy to recover past demographic history is based on the strategy described and used in [13, 75]. The idea of this approach is to generate a representative set of demographic scenarios over which the network generalizes to consequently infer similar demographic changes after training. More details on the training strategy can be found in Supplementary Text S2.

To improve the outputted demographic history, we introduce a smoothing of the demography allowing to infer continuous variation of population size through time. We do so by interpolating I time points cubically, and choosing w (set to 60) uniformly spaced new time points of the interpolation in log space. All time points more recent than ten generations in the past are discarded, since inference is too imprecise in the very recent present under our models. An example of this process can be seen in Supplementary Text S2.

Training dataset to disentangling coalescent and selection signatures

Beyond parameter inference, deep learning approaches can also be used for clustering. Hence, we train a GNN to disentangle between different scenarios and models. In total, we define eight classes, namely K (S0) (Kingman, no selection), K (WS) (Kingman, weak selection), K (MS) (Kingman, medium selection), K (SS) (Kingman, strong selection) and four different β -coalescent classes ([2.0-1.75], [1.75-1.5], [1.5-1.125], [1.125-1.01]) without selection. The three different selection regimes are defined, corresponding to $Ne \times s$ in: [0.1, 0.01] for SS, [0.01, 0.001] for MS, [0.001, 0.0001] for WS and [0] for absence of selection. Demography is kept constant and set to 10^5 individuals and sequence length is set to 10^5 bp. The simulation is discarded if it resulted in less than 2,000 obtained trees and is rerun with twice the sequence length until the tree number required is satisfied. This procedure avoids simulating large genome segments of which only a small fraction of trees is used for the given scenario during training and inference. The selection site is introduced in the centre of the respective sequence, so that 249 trees left and 250 right of the middle tree under selection form a training sample, using 500 trees for each sample. One hundred replicates are generated for each training sample. The complete training dataset consists of 1,000 parameter sets, 500 for the Kingman cases and 500 for the β -coalescent cases, with approximately 125 parameter sets per class. The model itself is trained on one epoch (number of time the data is analyzed), and the evaluation performed afterwards on 1,000 randomly generated parameter sets, with one replicate per parameter set. The same architecture used for demography estimation is employed with additional linear layers to reduce the number of output dimensions from 60 to 8. The loss function is set to a Cross-Entropy-Loss for the network to be trainable for categorical labels. Otherwise all architecture and training parameters is the same as described above and detailed in Supplementary Text S2.

Results

Inference bias under the wrongly assumed Kingman coalescent

We first study the effect of assuming a Kingman coalescent when the underlying true model is a β -coalescent (*i.e.* in presence of multiple merger events) by applying MSMC and MSMC2 to our simulated data. The inference results from MSMC and MSMC2 when the population undergoes a sawtooth demographic scenario are displayed in Figure 2. For $\alpha > 1.5$ the shape of the past demography is fairly well recovered. Decreasing the parameter α of the β -coalescent (*i.e.* higher probability of multiple merger events occurring) increases the variance of inferences and flattens the demography. Yet, both methods fail to infer the correct population size, due to the scaling discrepancy between the Kingman and β -coalescent. Hence, we perform the same analysis and correct for the scaling effect of the MMC versus a Kingman coalescent to better capture the specific effects of assuming binary mergers only. The results are displayed in Figure S1. For $\alpha > 1.5$ the demography is accurately recovered providing we know the true value of α to adjust the y-axis (population size) scale. However, for smaller α values the observed variance is extremely high and a flattened past variation of population size is observed.

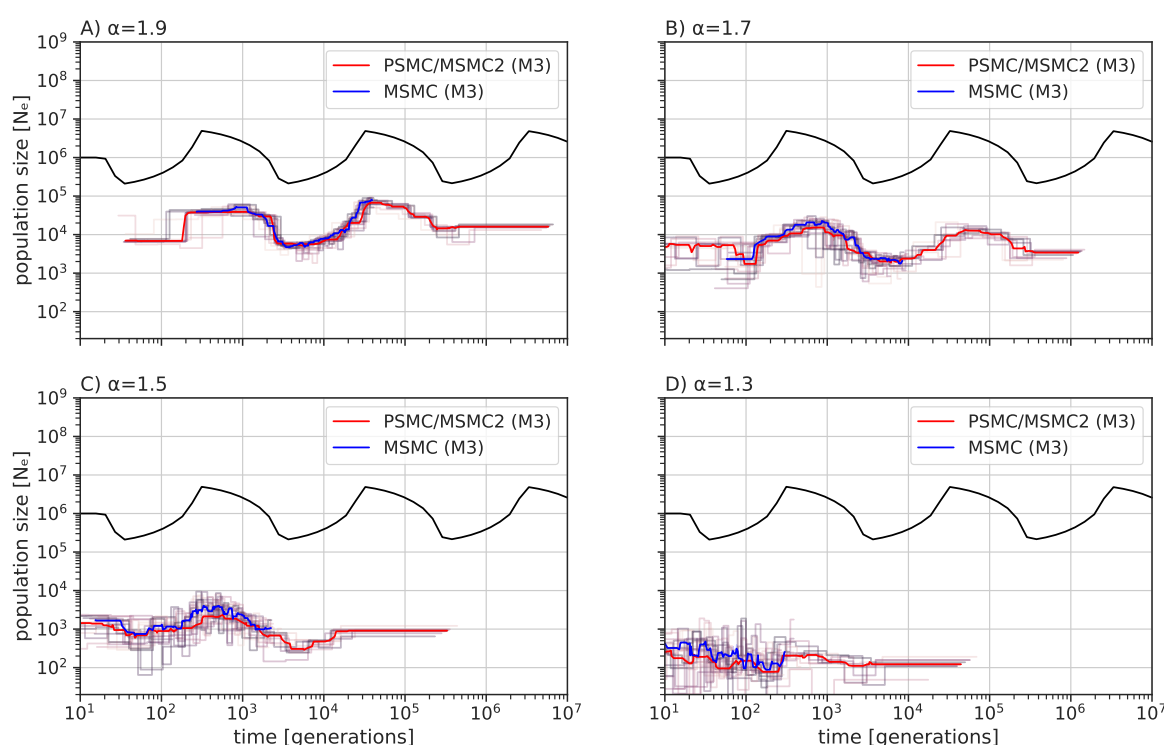


Fig. 2. Performance of MSMC and MSMC2 under a β -coalescent. Averaged estimated demographic history by MSMC (blue) and MSMC2 (red) based on 10 sequences of 100 Mb with $\mu = r = 10^{-8}$ per generation per bp over ten repetitions (while analyzing simultaneously 3 sequences, noted by M=3). Each repetition result is represented in light red (PSMC/MSMC2) or in light blue (MSMC). Population undergoes a sawtooth demographic scenario (black) for A) $\alpha = 1.9$, B) $\alpha = 1.7$, C) $\alpha = 1.5$, and D) $\alpha = 1.3$.

The limit of the Markovian hypothesis

As SMC approaches rely on the hypothesis of Markovian change in genealogy along the genome, we study the effect of α on the linkage disequilibrium (LD) of pairs of SNPs (r^2 , [71, 60]) in data simulated under the Kingman Coalescent or the β -coalescent (with $\alpha = 1.5$ and $\alpha = 1.3$) and constant population size (Figure 3). Linkage monotonously decreases with distance under the Kingman coalescent. Under the β -coalescent a similar shape of the distribution is observed but with a higher average amount of LD. We find a higher variance in LD for smaller α values. The increased variance results in the occurrence of high spikes of LD along the genome (*e.g.* Figure 3 B). The stochastic increase of linkage along the genome demonstrates that the Markovian hypothesis used to model genealogies along the genome is violated under the β -coalescent due to the long range effect of strong multiple merger events [8].

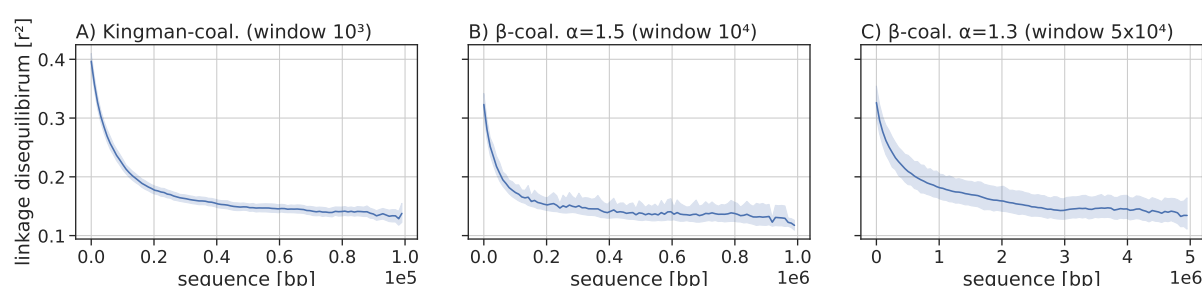


Fig. 3. Linkage disequilibrium under a Kingman and β -coalescent. Pairwise linkage disequilibrium between SNPs (r^2) under a Kingman and β -coalescent with $\alpha = 1.5$ and $\alpha = 1.3$ using 50 sequences of length A) 0.1 Mb, B) 1Mb, and C) 5Mb. The population size is constant at $N = 10^4$ for the Kingman model and $N = 10^6$ for the β -coalescent, with $\mu = 2 \times 10^{-7}$ and $r = 2 \times 10^{-8}$ per generation per bp. For each LD analysis, the linkage disequilibrium is calculated by averaging it over a window of size 10^3 , 10^4 and 5×10^4 bp respectively in A), B) and C).

We further investigate the effect of multiple merger events on LD. To this aim, we first assume an SMC framework (*e.g.* MSMC2 or eSMC) to predict the transition matrix (*i.e.* matrix containing the probabilities for the coalescent time to change to another value along the genome) and investigate the absolute difference between the observed transition events. Under the Kingman coalescent, the distribution of coalescent times between two positions in a sample of size two ($n = 2$) is well approximated by the SMC as shown in Figure S2 (*i.e.* absence of structured difference between observed and predicted). However, under the β -coalescent (with $\alpha = 1.3$) we observe significant differences between observed and predicted at times points where multiple merger events occur (Figure S3). In practice, multiple merger events do not occur at each time point (as they remain rare events), unveiling a discrepancy between the expectation from the SMC (*i.e.* approximating the distribution of genealogies along the genome by a Markov chain) and the simulated data. This discrepancy does not stem from the simulator, because it correctly generates ARG under the β -coalescent model [8, 7], but from the limits of the SMC approximation to model events with long range effects on the ARG (Figure S3).

Inferring α and past demography on ARG

To test if our two approaches (GNNcoal and SM β C) can recover the past variation of population size and the α parameter, we run both methods on simulated tree sequences under different α values and demographic scenarios. Figure 4 displays results for data simulated under a sawtooth past demography and for α ranging from 1.9, 1.7, 1.5 to 1.3. In all cases, the GNNcoal approach exhibits high accuracy and low variance to infer the variation of population size. For high α values (>1.5), the shape of population size variation is well recovered by SM β C. However, for smaller values, the observed high variance demonstrates the limits of SMC inferences. On average, both approaches seem to recover fairly well the true α value (Figure 4 and Table 1 in S1). In particular, GNNcoal displays high accuracy and lower standard deviation. We note that the variance in the estimation of α increases with diminishing α value. Moreover, increasing the number of simultaneously analyzed sequences by SM β C does not seem to improve the inferred α value (Table S1). These conclusions are also valid for the results in Figures S4-S7 and Table S1 based on inference under four additional demographic scenarios: constant population size, bottleneck, sudden increase and sudden decrease of population size.

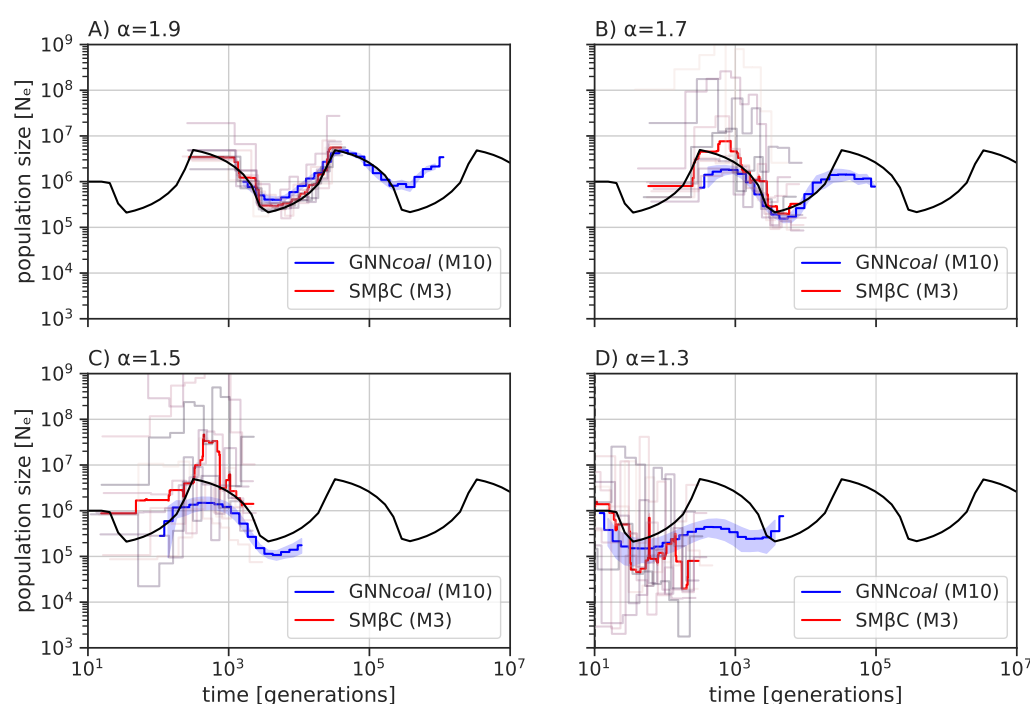


Fig. 4. Best-case convergence estimations of SM β C and GNNcoal under a β -coalescent. Estimations of past demographic history by SM β C in red (median) and by GNNcoal in blue (mean and 95% confidence interval, CI95; while analyzing simultaneously 3 or 10 sequences, noted by M=3 or M=10) when population undergoes a sawtooth demographic scenario (black) under A) $\alpha = 1.9$, B) $\alpha = 1.7$, C) $\alpha = 1.5$ and D) $\alpha = 1.3$. SM β C runs on 10 sequences and 100 Mb, GNNcoal runs on 10 sequences and 500 trees, and $\mu = r = 10^{-8}$ per generation per bp.

Because when α diminishes, the effective population size decreases and the number of recombination events plummets for small values of $\alpha < 1.5$. To demonstrate the

Inference under the Beta Coalescent

theoretical convergence of $SM\beta C$ to the correct values, we run $SM\beta C$ on data simulated with mutation and recombination rate fifty times higher under similar scenarios as in Figure 4. This operation increases the amount of data in the form of SNPs and number of independent coalescent trees by recombination. Results of $SM\beta C$ for α values of 1.7, 1.5 and 1.3 are displayed on Table S2. Results show that $SM\beta C$ can recover α with higher accuracy when more data is available.

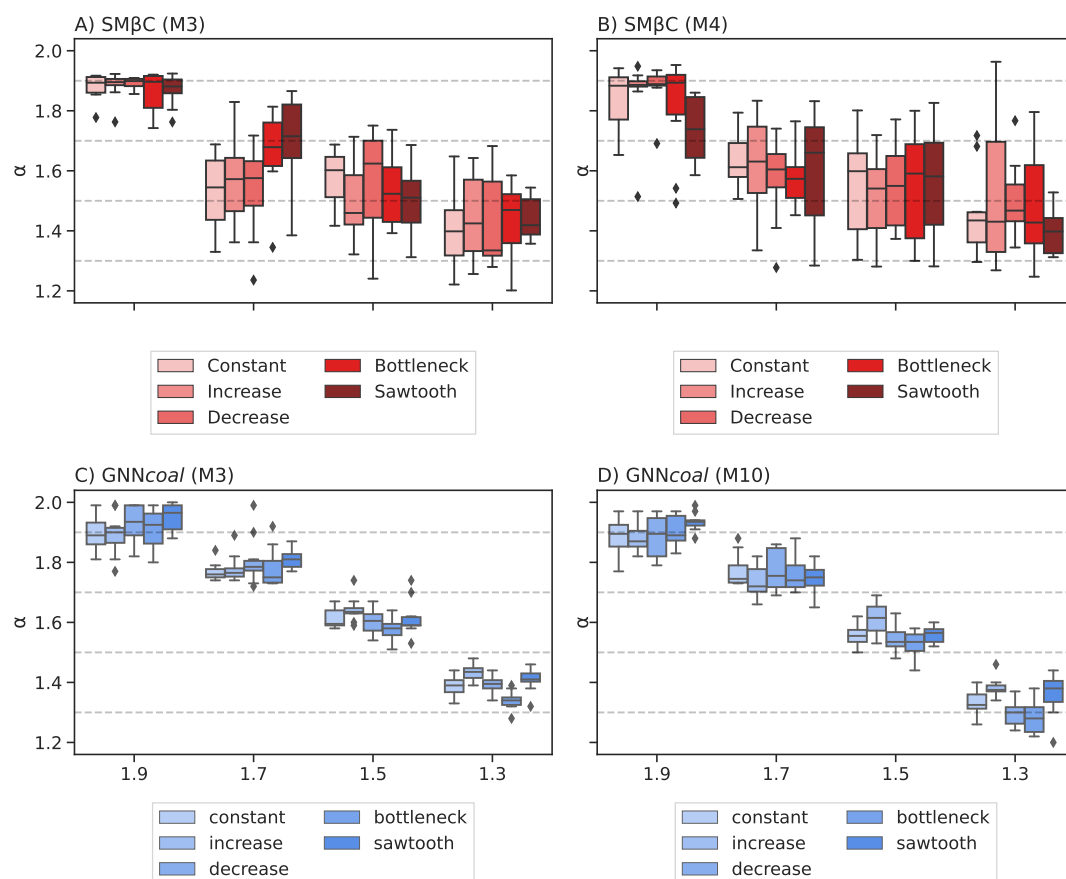


Fig. 5. Estimated α values by $SM\beta C$ and $GNNcoal$. Estimated values of α by $SM\beta C$ and $GNNcoal$ over ten repetitions using 10 sequences of 100 Mb with $\mu = r = 10^{-8}$ per generation per bp under a β -coalescent process (with different α parameter). The analysis are run on five different demographic scenarios (Constant population size, Bottleneck, Sudden increase, Sudden decrease and a Sawtooth demography) using a sample size $n = 3$ for A) and C), $n = 4$ for B), and $n = 10$ for D). Grey dashed lines indicate the true α values.

Although 10 sequences are given to $SM\beta C$ in the previous analyses, the method can only analyze three or four simultaneously. On the other hand, $GNNcoal$ can simultaneously analyze 10 sequences, that is the whole simulated ARG. As we observe that $GNNcoal$ has a higher performance than $SM\beta C$, we wish to test whether the $GNNcoal$ better leverages information from the ARG or benefits from simultaneously analyzing a larger sample size. Thus, we run $GNNcoal$ on the same dataset, but downsampling the coalescent trees to a sample size three. Results are displayed in Figure S4 to Figure S7. Results with sample size three of $GNNcoal$ are similar to results with sample size 10,

Inference under the Beta Coalescent

demonstrating that the GNNs can better leverage information from the ARG in presence of multiple merger events (Figure S8).

Additionally, we test if both approaches can recover a Kingman coalescent from the ARG when data are simulated under the Kingman coalescent, namely both approach should recover $\alpha = 2$. To do so, we simulate the same five demographic scenarios as above under a Kingman coalescent and infer the α parameter along with the past variation of population size. Estimations of α values are provided in Table 1 and are systematically higher than 1.85, suggesting mostly binary mergers. The associated inferred demographies are shown in Figures S9-S13. Both approaches correctly infer the past demographic shape up to the scaling discrepancy between the Beta and the Kingman coalescent (as previously described). Furthermore, we notice that the scaling effect only affects the y-axis for the SM β C but affect both axes for GNN $coal$.

As GNN $coal$ was not trained on data simulated under the Kingman coalescent (especially with such high population size), some events fall beyond the scope of the GNN due to the scaling discrepancy between the Beta and Kingman coalescence. Hence, we run GNN $coal$ on data simulated under the Kingman coalescent but with smaller population size (scaled down by a factor 100) to assure that all events fall within the scope of the GNN. Values of α inferred by the GNN $coal$ and the SM β C under the five demographic scenarios are available in Table S3. The associated inference of population size are plotted in Figure S9-S12. Both approaches recover high α values (*i.e.* > 1.85) suggesting a genealogy with almost exclusively binary mergers. In addition, both approaches accurately recover the shape of the past variation of population size up to a scaling constant but only on the population size y-axis.

Inferring α and past demography from simulated sequence data

We first investigate results for both GNN $coal$ and SM β C when the ARG is reconstructed with ARGweaver [69], the latter being considered the best performing approach to infer ARG for sample size smaller than 20 [15]. Demographic inference results by both approaches are displayed in Figure S14, and α inference results in Table S4. GNN $coal$ does not recover the shape of the demographic history from the inferred ARGs and largely overestimates α . In contrast, SM β C produces better inferences of α when giving the inferred ARG as input. SM β C recovers the shape of the past variation of population size for $\alpha > 1.3$ but displays extremely high variance for $\alpha = 1.3$.

We then run SM β C on simulated sequence data and found that α is typically well recovered (Table 2) and that results are similar to what obtained when the true ARG is given. Furthermore, the shape of the past variation of population size is well inferred under the sawtooth demographic scenario for $\alpha > 1.3$ (Figure S15). In the other four scenarios, the shape of the demography is recovered in recent times but population sizes are underestimated in the past (Figure S16). Finally, as found above from inputted ARGs, the variance in estimates of population sizes generally increases with diminishing α .

Inferring MMC and accounting for selection

As specific reproductive mechanisms and selection can lead to the occurrence of multiple merger-like events, we train our neural network on data simulated under the β -coalescent, and under the Kingman coalescent in presence or absence of selection. We then use the trained *GNNcoal* to determine if multiple merger events originate from skewed offspring distribution or positive selection, or if the data follows a neutral Kingman coalescent process. The classification results are displayed in Figure 6 in the form of a confusion matrix, that is the percentage of times the *GNNcoal* correctly assigns the true model shown on the diagonal evaluated on a test dataset of 1,000 known ARGs. Our approach can accurately select the model except in two cases. *GNNcoal* shows limited power to distinguish between strong and intermediate selection under the Kingman coalescent, as well as to distinguish between the β -coalescent with a small amount of multiple merger events (*i.e.* $\alpha > 1.75$) and the β -coalescent case with $1.75 > \alpha > 1.5$.

Inference under the Beta Coalescent

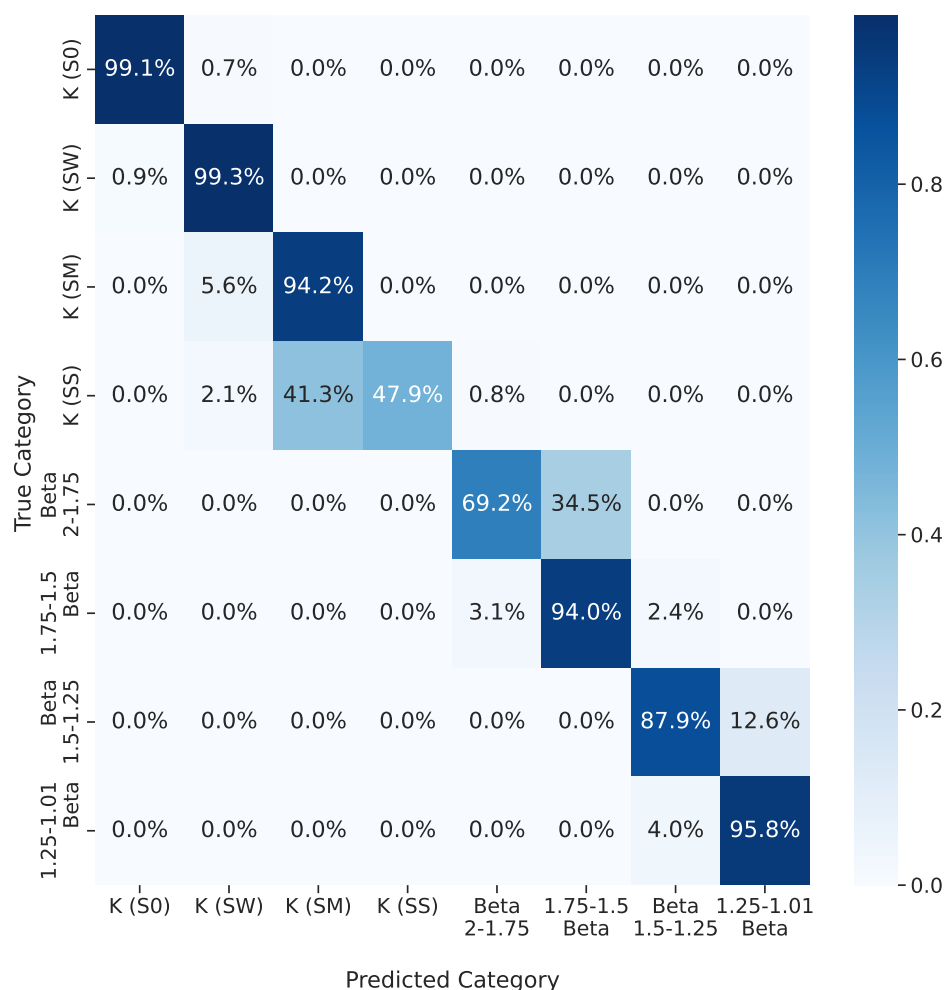


Fig. 6. Confusion matrix for Kingman and β -coalescent classification model under varying selection coefficients. Evaluation of classification accuracy for Kingman (K) and β -coalescent (B) for no selection (S0), weak selection (SW), medium selection (SM) and strong selection (SS) using a 1,000 repetition validation dataset. Population size was kept constant at $N = 10^5$ individuals, using a sample size $n = 10$ and $r = 10^{-8}$ per bp per generation.

To assess the effect of selection, we infer α along the genome with both approaches from data simulated with strong positive selection or neutrality under a Kingman coalescent with population size being constant through time. SM β C infers α on windows of 10kbp along the genome, and GNNcoal infers α every 20 trees along the genome. Results for GNNcoal and SM β C are displayed in Figure 7. Both approaches recover smaller α value around the locus under strong selection. However under neutrality or weak selection, inferred α values remain high (>1.6) along the genome.

Inference under the Beta Coalescent

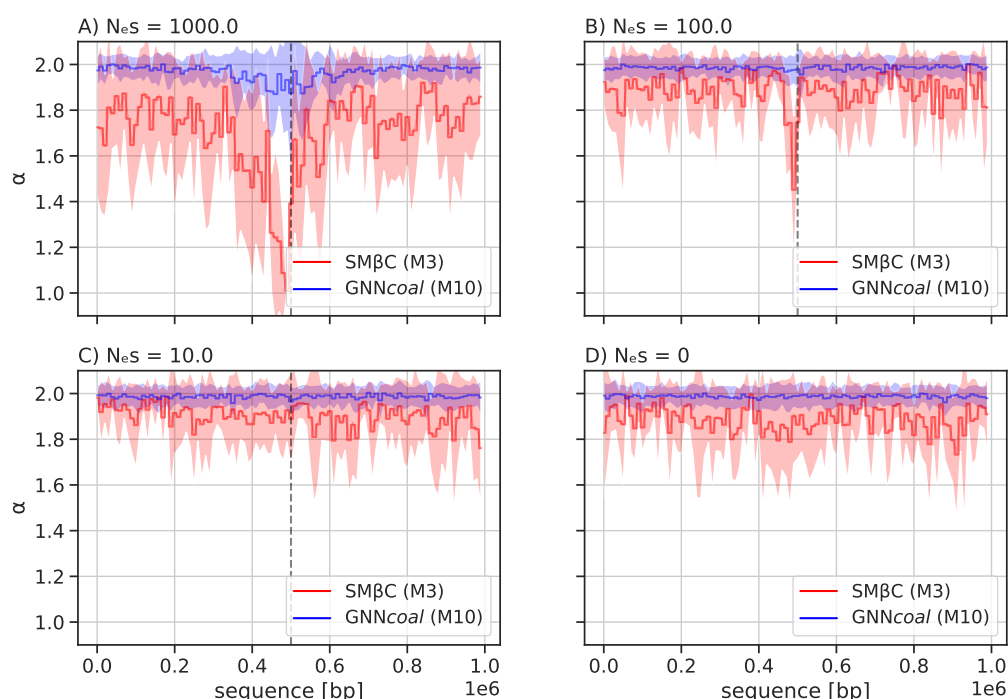


Fig. 7. Averaged estimations by GNNcoal and SM β C under selection Estimations of α along the genome by the GNNcoal approach and the SM β C when population undergoes as strong positive selective sweep event (at position 0.5 Mb) under different strengths of selection: A) $s = 0.01$, B) $s = 0.001$, C) $s = 0.0001$, and D) $s = 0$ meaning neutrality (mean and standard deviation for both methods). The population size is constant and set to $N = 10^5$ with $\mu = r = 10^{-8}$ per generation per bp. We hence have in A) $N_e \times s = 1000$, B) $N_e \times s = 100$, C) $N_e \times s = 10$ and D) $N_e \times s = 0$. SM β C uses 20 sequences of 1Mb (red) and GNNcoal uses 10 sequences through down-sampling the sample nodes (blue)

Similarly, we run both approaches on data simulated under the β -coalescent (assuming neutrality) and we infer the α value along the genome. Inferred α values by both approaches are plotted in Figure 17 in S1. GNNcoal is able to recover the α value along the genome with moderate overestimation due to tree sparsity. On the contrary, SM β C systematically underestimates α values. Nevertheless, unlike in presence of positive selection at a given locus, the inferred α values are found in all cases to be fairly constant along the genome.

We finally simulate data under a strong selective sweeps or under neutrality conditioned on a sawtooth demographic scenario. Under neutrality, our both approaches recover, as expected, high α values along the genome and can accurately recover the past variation of population size (only up to a scaling constant for GNNcoal, since it was trained on the β -coalescent only) (Figure 8). Similarly, when the simulated data contains strong selection, a small α value is recovered at the locus under selection and the past variation of population size is accurately recovered, albeit with a small underestimation of population size in recent times for SM β C (Figure 8).

Inference under the Beta Coalescent

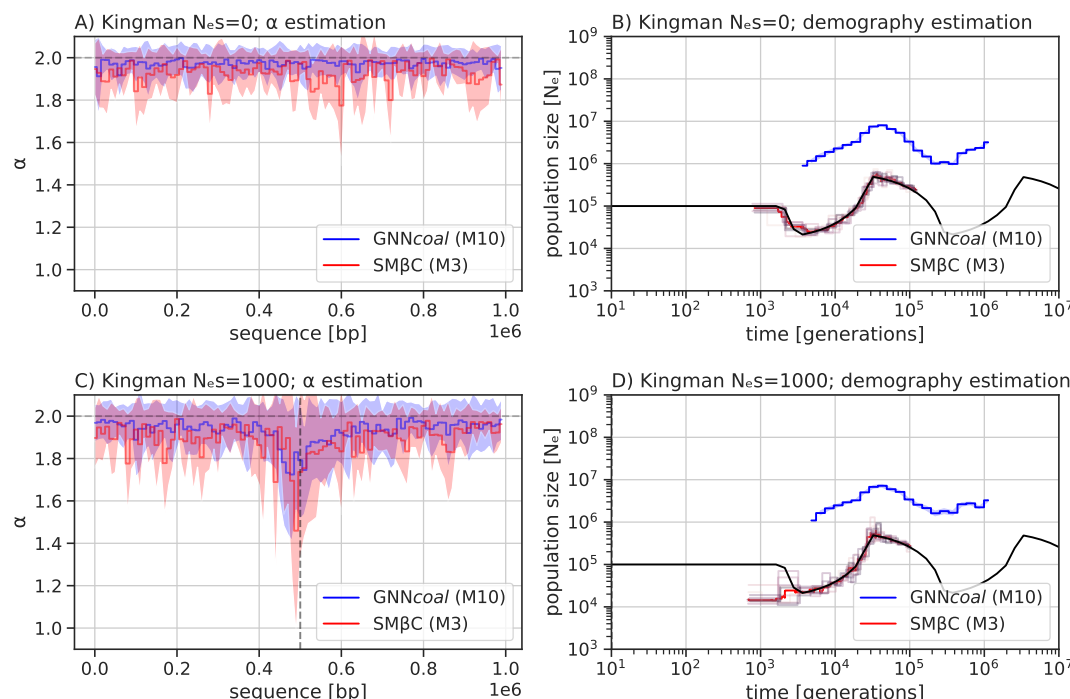


Fig. 8. Simultaneous estimations of α along the sequence under demographic change by GNNcoal and SM β C. Simultaneous estimation of α along the genome under a partial sawtooth scenario: A) and B) in the absence of selection (mean and standard deviation for both methods), and C) and D) presence of selection with $N_e S = 1,000$ (mean and CI95 for GNNcoal and median for SM β C). SM β C uses 20 sequences of 1Mb (red) and GNNcoal uses 10 sequences through down-sampling the sample nodes (blue), and $\mu = r = 10^{-8}$ per generation per bp.

Discussion

With the rise in popularity of SMC approaches for demographic inferences [54], most current methods leverage information from whole genome sequences by simultaneously reconstructing a portion of the ARG to infer past demographic history [54, 78, 89, 90], migration rates [47, 91], variation in recombination and mutation along the genome [5, 4], as well as ecological life history traits such as selfing or seed banking [81, 87]. However, other previous studies proposed to uncouple both steps, namely by first reconstructing the ARG and by then inferring parameters from its distribution [82, 32, 69]. Indeed, recent efforts have been made to improve approaches to recover the ARG [84, 46, 37, 69, 55, 15], as well as its interpretation [31, 82]. Our results on data simulated under the β -coalescent clearly show the strong effect of multiple merger events on the topology and branch length of the ARG. We find that the more multiple merger events occur, the more information concerning the past demography is lost. Both GNNcoal and SM β C, whether given sequence data, the true or inferred ARG, can recover the α parameter and the variation of past population size for α values high enough (*i.e.* $\alpha > 1.3$). However, for lower values of α , a larger amount of data is necessary for any inference which becomes nearly impossible when α tends to one. Both approaches can also recover the Kingman coalescent (*i.e.* $\alpha > 1.8$). We find that GNNcoal outperforms SM β C in almost all cases when given the true ARG, and we demonstrate that GNNcoal can be used to disentangle

Inference under the Beta Coalescent

between β -coalescent and Kingman models with selection.

Overall, our results provide a substantial improvement in the development of inference methods for models with multiple merger events, a key step to understand the underlying reproduction mechanism of a species. While still inferring population sizes of the correct order of magnitude, $SM\beta C$ is outperformed by $GNNcoal$ when given true ARGs as input. As we directly compare our theoretical SMC to the GNN based on the same input data (coalescent trees), we are ideally placed to dissect the mechanisms underlying the power of the $GNNcoal$ method. We identify four main reasons for the difference in accuracy between the two methods developed. First, the $SM\beta C$ approach suffers from the limit of the sequential Markovian coalescent hypothesis along the genome when dealing with strong multiple merger events [8, 21]. Second, most of current SMC approaches rely on a discretization of the coalescent times into hidden states, meaning that simultaneous mergers of three lineages may not be easily distinguished from two consecutive binary mergers occurring over a short period of time. Third, the $SM\beta C$ relies on a complex hidden Markov model and due to computational and mathematical tractability, it cannot leverage information on a whole genealogy. In fact, as MSMC, $SM\beta C$ only focuses on the first coalescent event, and therefore cannot simultaneously analyze large sample size. Furthermore, the $SM\beta C$ approach leverages information from the distribution of genealogies along the genome. Whilst, in the near absence of recombination events, both approaches cannot utilize any information from the genealogy itself, $GNNcoal$ can overcome this limit by increasing the sample size. Fourth, the $SM\beta C$ is based on a coalescent model where α is constant in time. Yet multiple merger events do not appear regularly across the genealogical timescale, but occur at few random time points. Hence, the SMC approach suffers from a strong identifiability problem between the variation of population size and the α parameter (for low α values). For instance, if during one hidden state one strong multiple merger event occurs, multiple merger events are seldom observed and $SM\beta C$ may rather assume a small population size at this time point (hidden state). This may explain the high variance of inferred population sizes under the β -coalescent.

By contrast, $GNNcoal$ makes use of the whole ARG, and can easily scale to larger sample sizes (over 10), although it recovers α with high accuracy with sample size 3 only. Our interpretation is that $GNNcoal$ is able of simultaneously leveraging information from topology and the age of coalescent events (nodes) across several genealogies (here 500). $GNNcoal$ ultimately leverages information from observing recurrent occurrences of the same multiple merger events at different locations on the genome, while being aware of true multiple merger events from rapid successive binary mergers. We believe that our results pave the way towards the interpretability of GNN and deep learning methods applied to population genetics.

When applying both approaches to simulated sequence data (and not to true ARGs), both approaches behave differently. $GNNcoal$ is not capable to accurately infer model parameters, i.e. past variation of population size or α . In contrast, $SM\beta C$ performed better than $GNNcoal$ when dealing with sequence data (and not true ARG). $SM\beta C$ is capable of recovering α and the shape of the demographic scenario in recent times irrespective of whether sequence data or ARG inferred by ARGweaver is given as input. This is most likely because the statistic used by $SM\beta C$ (i.e. first coalescent event in discrete time) is coarser than the statistic used by $GNNcoal$ (i.e. the exact ARG). We therefore speculate

Inference under the Beta Coalescent

that the theoretical framework of the $SM\beta C$, although being in theory less accurate than *GNNcoal*, is more robust and suited for application to sequence data. More specifically, the issue being faced by the *GNNcoal* is known as out-of-distribution inference [39], which requires the network to generalize over an untrained data distribution. This issue happens because *GNNcoal* is not trained using ARG inferred by ARGweaver. Building a training data set for *GNNcoal* to overcome this issue is currently impractical due to the inference speed of ARGweaver. However, future work will aim at increasing robustness of GNN inferences, for instance by adding uncertainty or multiple models during the training process. Improving the performance of *GNNcoal* on sequence data requires more efficient and accurate ARG inference methods (, which can be used on a broader spectrum of data sets. The latter observation is important to avoid bias from potential hypothesis violations of the chosen ARG inference approach).

Past demographic history, reproductive mechanisms, and natural selection are among the major forces driving genome evolution [41]. Hence, in the second part of this manuscript we focus on integrating selection in both approaches. Currently, no method (especially if relying only on SFS information) can account for the presence of selection, linkage, non-constant population size and multiple merger events [41] although recent theoretical framework might render this possible in the future [1]. As a first step to fill this gap, we demonstrate that *GNNcoal* can be used for model selection to reduce the number of hypotheses to test. Determining which evolutionary forces are driving the genome evolution is key, as only under the appropriate neutral population model results of past demography and selection scans can be correctly interpreted [41, 43]. The high accuracy of *GNNcoal* in model selection is promising, especially as other methods based on the SFS alone [52, 44] have limits in presence of complex demographic scenarios. GNN can possibly overcome these limits, as it is easier to scale the GNN to estimate more parameters. We follow a thread of previous work [72, 36, 11], by integrating and recovering selection, multiple merger and population size variation by simply allowing each fixed region in the genome to have its own α parameter. In presence of strong selection, we find lower α value around the selected loci and high α value in neutral neighbouring regions. In presence of weak selection, no effect on the estimated α value is observed, demonstrating that weak selection can be modeled by a binary merger and has only a local effect on the branch length by shortening it. Hence, our results point out that strong selection can indeed be modeled as a local multiple merger event (see [24, 11, 72]). In theory, both approaches should be able to infer the global α parameter linked to the reproductive mechanism, as well as the local α parameter resulting from selection jointly with the variation of population size. However, the absence of a simulator capable of simulating data with selection and non-constant population size under a β -coalescent model prevents us from delivering such proofs. We show strong evidence that under neutrality our approaches can recover a constant (and correct) α along the genome as well as the past variation of the population size. We further predict that, while selective processes favor coding regions, local variations in α (as a consequence of sweepstake events) should be indifferent to coding or non-coding regions. Hence, we suggest that current sequence simulators [7, 33] could be extended to include the aforementioned factors and *de facto* facilitate the development of machine learning approaches.

Our study is unique in developing a new state-of-the-art SMC approach and demonstrating that computational and mathematical problems can be overcome by deep learn-

Inference under the Beta Coalescent

ing (here GNN) approaches. The *GNNcoal* approach is, in principle, not limited to the β -coalescent, and should work for other multiple merger models (*e.g.*, Dirac coalescents [25]) with the appropriate training. Furthermore, our *SM β C* approach is the first step to build a full genome method with an underlying model accounting for positive selection. In the future, further implementations may be added for a more realistic approach. The α parameter should be varying along the genome (as a hidden state), as the recombination rate in the iSMC [5]. This would allow to account for the local effect of strong and weak selection [1]. The effect of the α parameter could be also changing through time to better model the non uniform occurrence of multiple merger events through time. Although it is mathematically correct to have α as a constant in time, it is erroneous in practice (Figure 2 in S1). We speculate that those additional features will allow to accurately model and infer multiple merger events, variation of population size, and selection at each position on the genome. We believe that deep learning approaches could also be improved to recover more complex scenarios, providing in depth development on the structure of the graph neural networks, for example, by accounting for more features. At last, further investigation are required to make progress in the interpretability of the GNN methods, namely which statistics and convolution of statistics are used by *GNNcoal* to infer which parameters.

As our approaches are the first of their kind, we chose to restrain our study to haploid models of β and Kingman coalescent as a proof of principle. However, the *GNNcoal* and *SM β C* approaches can be extended to higher ploidy levels. Diploid versions of the haploid reproduction models whose genealogies are given by the β -coalescent lead to slightly different MMC coalescent models which can exhibit simultaneous multiple mergers [8, 10]. Thus, our GNN approach should be directly applicable when trained on these diploid models which are implemented in *msprime* [7]. However, to adjust the *SM β C* approach would be somewhat more cumbersome (but doable), since we would need to extend the underlying HMM to account for simultaneous multiple mergers. We emphasise that while there is growing evidence that MMC models produce better fitting genealogies for various species [30], there is ongoing discussions about which mathematical models are better suited to which species (for example see [3] for cod). We advocate that the life-cycle and various ecological factors determine whether a haploid or diploid MMC model can be chosen. On the one hand, a diploid MMC model is likely realistic if the species has a diploid life-cycle and balanced sex-ratio, so that multiple merger events do indeed happen in both sexes. On the other hand, if species are mostly haploid or clonal/asexual during their life-cycle (with periodically one short diploid phase for sexual reproduction) or exhibit strongly imbalanced sex-ratio, a haploid MMC model may be better suited. In their current form, our approaches are applicable to data from species with the latter characteristics such as many fungal and micro-parasites of plants and animals (including humans) as well as invertebrates (*e.g.* *Daphnia* or aphids) which undergo several clonal or parthenogenetic phases of reproduction (and one short sexual phase) per year. This represents a non-negligible set of study organisms which are of importance for medicine and agriculture [88].

Our results on inferred ARGs stress the need for improving ARG inference [15]. Thanks to the SMC we are close to model the ARG allowing to infer demographic history, selection and specific reproductive mechanism. Moreover, the comparison of deep learning approaches with model driven *ad hoc* SMC methods may have the potential to help us solve ongoing challenges in the field. These include simultaneously inferring and

Inference under the Beta Coalescent

accounting for recombination, variation of population size, different type of selection, population structure and the variation of the mutation and recombination rate along the genome. These issues have puzzled theoreticians and statisticians since the dawn of population genetics [41].

On a final note, as environmental changes hit us all, we suggest that decreasing the computer and power resources needed to perform DL/ GNN analyses should be attempted [76]. Based on our study, we suggest that population genetics DL methods could be built as a two step process: 1) inferring ARGs, and 2) inferring demography and selection based on the ARGs. We speculate that general training sets based on ARGs could be build and be widely applicable for inference across many species with different life cycles and life history traits, while the inference of ARGs could be undertaken by complementary deep learning or Hidden Markov methods.

Tables

| scenario | True α | α :SM β C,M=3 | α :SM β C,M=4 | α : GNN, M=3 | α : GNN, M=10 |
|------------|---------------|----------------------------|----------------------------|---------------------|----------------------|
| Constant | 2 | 1.97 (0.005) | 1.97 (0.008) | 1.99 (0.002) | 1.99 (0.003) |
| Sawtooth | 2 | 1.94 (0.017) | 1.87 (0.019) | 1.99 (0.002) | 1.99 (0.003) |
| Bottleneck | 2 | 1.97 (0.01) | 1.97 (0.009) | 1.99 (0.003) | 1.99 (0.004) |
| Decrease | 2 | 1.97 (0.007) | 1.97 (0.008) | 1.99 (0.003) | 1.99 (0.004) |
| Increase | 2 | 1.97 (0.007) | 1.97 (0.008) | 1.99 (0.004) | 1.99 (0.002) |

Table 1: Average estimated values of α by SM β C and GNN $coal$ over ten repetitions under the Kingman coalescent using 10 haploid sequences of 10 Mb and $\mu = r = 10^{-8}$ per generation per bp. The standard deviation is indicated in brackets.

| scenario | True α | α^* :SM β C,M=3 |
|------------|---------------|------------------------------|
| Constant | 1.9 | 1.86 (0.16) |
| Bottleneck | 1.9 | 1.89 (0.09) |
| Increase | 1.9 | 1.93 (0.07) |
| Decrease | 1.9 | 1.96 (0.04) |
| Sawtooth | 1.9 | 1.76 (0.17) |
| Constant | 1.7 | 1.82 (0.10) |
| Bottleneck | 1.7 | 1.64 (0.23) |
| Increase | 1.7 | 1.82 (0.10) |
| Decrease | 1.7 | 1.89 (0.13) |
| Sawtooth | 1.7 | 1.71 (0.27) |
| Constant | 1.5 | 1.52 (0.30) |
| Bottleneck | 1.5 | 1.64 (0.33) |
| Increase | 1.5 | 1.57 (0.24) |
| Decrease | 1.5 | 1.60 (0.18) |
| Sawtooth | 1.5 | 1.66 (0.14) |
| Constant | 1.3 | 1.31 (0.20) |
| Bottleneck | 1.3 | 1.2 (0.17) |
| Increase | 1.3 | 1.24 (0.13) |
| Decrease | 1.3 | 1.57 (0.11) |
| Sawtooth | 1.3 | 1.37 (0.16) |

Table 2: Average estimated α values by SM β C on simulated sequence data over ten repetitions using 10 sequences of 10 Mb with recombination and mutation rate set to 1×10^{-8} for α 1.9 and 1.7, 1×10^{-7} for α 1.5 and 1×10^{-6} for α 1.3 per generation per bp under a Beta coalescent process. The analysis are run on five different demographic scenarios (Constant population size, Bottleneck, Sudden increase, Sudden decrease and a Sawtooth demography).

Data availability

Code used to generate the simulated data for analysis, training and validation alongside (trained) deep learning models can be found at <https://github.com/kevinkorfmann/GNNcoal> and <https://github.com/kevinkorfmann/GNNcoal-analysis>. Code for SMC approaches used in this manuscript are available in the R package eSMC2 <https://github.com/TPPSellinger/eSMC2>.

Acknowledgments

This work was supported by the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure (de.NBI) (031A532B, 031A533A, 031A533B, 031A534A, 031A535A, 031A537A, 031A537B, 031A537C, 031A537D, 031A538A). KK is supported by a grant from the Deutsche Forschungsgemeinschaft (DFG) through the TUM International Graduate School of Science and Engineering (IGSSE), GSC 81, within the project GENOMIE QADOP. TS is supported by the Austrian Science Fund (project no. TAI 151-B). AT acknowledges funding from the DFG grant TE809/1-4 (project 254587930) and TE809/7-1 (project 317616126). FF and AT acknowledge funding from the DFG Priority Program SPP1590 on "Probabilistic Structures in Evolution". MF and AT acknowledge the support from the Imperial College - TUM Partnership award.

Competing interests

The authors declare that no competing interests exist.

References

- [1] Frederic Alberti, Carolin Herrmann, and Ellen Baake. Selection, recombination, and the ancestral initiation graph. *THEORETICAL POPULATION BIOLOGY*, 142:46–56, DEC 2021.
- [2] Einar Arnason and Katrin Halldorsdottir. Nucleotide variation and balancing selection at the Ckma gene in Atlantic cod: analysis with multiple merger coalescent models. *PEERJ*, 3, FEB 24 2015.
- [3] Einar Árnason, Jere Koskela, Katrín Halldórsdóttir, and Bjarki Eldon. Sweepstakes reproductive success via pervasive and recurrent selective sweeps. *Elife*, 12:e80781, 2023.
- [4] Gustavo V. Barroso and Julien Y. Dutheil. Mutation rate variation shapes genome-wide diversity in *Drosophila melanogaster*. preprint, Evolutionary Biology, September 2021.
- [5] Gustavo V. Barroso, Natasa Puzovic, and Julien Y. Dutheil. Inference of recombination maps from a single pair of genomes and its application to ancient samples. *PLOS Genetics*, 15(11), NOV 2019.
- [6] CJ Battey, Peter L Ralph, and Andrew D Kern. Predicting geographic location from genetic variation with deep neural networks. *eLife*, 9:e54507, June 2020.

- [7] Franz Baumdicker, Gertjan Bisschop, Daniel Goldstein, Graham Gower, Aaron P. Ragsdale, Georgia Tsambos, Sha Zhu, Bjarki Eldon, E. Castedo Ellerman, Jared G. Galloway, Ariella L. Gladstein, Gregor Gorjanc, Bing Guo, Ben Jeffery, Warren W. Kretzschmar, Konrad Lohse, Michael Matschiner, Dominic Nelson, Nathaniel S. Pope, Consuelo D. Quinto-Cortes, Murillo F. Rodrigues, Kumar Saunack, Thibaut Sellinger, Kevin Thornton, Hugo van Kemenade, Anthony W. Wohms, Yan Wong, Simon Gravel, Andrew D. Kern, Jere Koskela, Peter L. Ralph, and Jerome Kelleher. Efficient ancestry and mutation simulation with msprime 1.0. *GENETICS*, 220(3), MAR 3 2022.
- [8] Matthias Birkner, Jochen Blath, and Bjarki Eldon. An Ancestral Recombination Graph for Diploid Populations with Skewed Offspring Distribution. *Genetics*, 193(1):255–290, JAN 2013.
- [9] Matthias Birkner, Jochen Blath, Martin Moehle, Matthias Steinruecken, and Johanna Tams. A modified lookdown construction for the Xi-Fleming-Viot process with mutation and populations with recurrent bottlenecks. *arXiv:0808.0412*, 2008.
- [10] Matthias Birkner, Huili Liu, and Anja Sturm. Coalescent results for diploid exchangeable population models I. *Electronic Journal of Probability*, 23, 2018.
- [11] Gertjan Bisschop, Konrad Lohse, and Derek Setter. Sweeps in time: leveraging the joint distribution of branch lengths. *GENETICS*, 219(2), OCT 2021.
- [12] Jochen Blath, Adrian Gonzalez Casanova, Noemi Kurt, and Maite Wilke-Berenguer. The seed bank coalescent with simultaneous switching. *Electronic Journal of Probability*, 25, 2020.
- [13] Simon Boitard, Willy Rodríguez, Flora Jay, Stefano Mona, and Frédéric Austerlitz. Inferring population size history from large samples of genome-wide molecular data - an approximate bayesian computation approach. 12(3):e1005877.
- [14] Erwin Bolthausen and A-S Sznitman. On ruelle’s probability cascades and an abstract cavity method. *Communications in mathematical physics*, 197(2):247–276, 1998.
- [15] Debora Y. C. Brandt, Xinzhu Wei, Yun Deng, Andrew H. Vaughn, and Rasmus Nielsen. Evaluation of methods for estimating coalescence times using ancestral recombination graphs. *GENETICS*, 221(1), MAY 5 2022.
- [16] Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, jul 2017.
- [17] E. Brunet, B. Derrida, A. H. Mueller, and S. Munier. Noisy traveling waves: Effect of selection on genealogies. *Europhysics Letters*, 76(1):1–7, OCT 2006.
- [18] E. Brunet, B. Derrida, A. H. Mueller, and S. Munier. Effect of selection on ancestry: An exactly soluble case and its phenomenological generalization. *Physical Review E*, 76(4, 1), OCT 2007.

- [19] Klara Elisabeth Burger, Peter Pfaffelhuber, and Franz Baumdicker. Neural networks for self-adjusting mutation rate estimation when the recombination rate is unknown. *PLOS Computational Biology*, 18(8):1–17, 08 2022.
- [20] Wenming Cao, Zhiyue Yan, Zhiquan He, and Zhihai He. A comprehensive survey on geometric deep learning. *IEEE Access*, 8:35929–35949, 2020.
- [21] Adrián González Casanova, Verónica Miró Pina, and Arno Siri-Jégousse. The Symmetric Coalescent and Wright-Fisher models with bottlenecks. *arXiv:1903.05642 [math]*, September 2020. arXiv: 1903.05642.
- [22] Jianhai Chen, Pan Ni, Xinyun Li, Jianlin Han, Ivan Jakovlic, Chengjun Zhang, and Shuhong Zhao. Population size may shape the accumulation of functional mutations following domestication. *BMC Evolutionary Biology*, 18, JAN 19 2018.
- [23] P Donnelly and TG Kurtz. Particle representations for measure-valued population models. *Annals of Probability*, 27(1):166–205, JAN 1999.
- [24] R Durrett and J Schweinsberg. A coalescent model for the effect of advantageous mutations on the genealogy of a population. *Stochastic Processes and their Applications*, 115(10):1628–1657, OCT 2005.
- [25] B Eldon and J Wakeley. Coalescent processes when the distribution of offspring number among individuals is highly skewed. *Genetics*, 172(4):2621–2633, APR 2006.
- [26] Bjarki Eldon, Matthias Birkner, Jochen Blath, and Fabian Freund. Can the Site-Frequency Spectrum Distinguish Exponential Population Growth from Multiple-Merger Coalescents? *Genetics*, 199(3):841+, MAR 2015.
- [27] Malaspinas et al. A genomic history of Aboriginal Australia. *Nature*, 538(7624):207+, OCT 13 2016.
- [28] Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with PyTorch geometric.
- [29] Lex Flagel, Yaniv Brandvain, and Daniel R Schrider. The Unreasonable Effectiveness of Convolutional Neural Networks in Population Genetic Inference. *Molecular Biology and Evolution*, 36(2):220–238, 12 2018.
- [30] Fabian Freund, Elise Kerdoncuff, Sebastian Matuszewski, Marguerite Lapierre, Marcel Hildebrandt, Jeffrey D. Jensen, Luca Ferretti, Amaury Lambert, Timothy B. Sackton, and Guillaume Achaz. Interpreting the pervasive observation of U-shaped Site Frequency Spectra. preprint, *Evolutionary Biology*, April 2022.
- [31] L. M. Gattepaille, M. Jakobsson, and M. G. B. Blum. Inferring population size changes with sequence and SNP data: lessons from human bottlenecks. *Heredity*, 110(5):409–419, MAY 2013.
- [32] Lucie Gattepaille, Torsten Guenther, and Mattias Jakobsson. Inferring Past Effective Population Size from Distributions of Coalescent Times. *Molecular Biology and Evolution*, 204(3):1191+, NOV 2016.

- [33] Benjamin C. Haller, Jared Galloway, Jerome Kelleher, Philipp W. Messer, and Peter L. Ralph. Tree-sequence recording in slim opens new horizons for forward-time simulation of whole genomes. *MOLECULAR ECOLOGY RESOURCES*, 19(2):552–566, MAR 2019.
- [34] Rebecca B. Harris and Jeffrey D. Jensen. Considering genomic scans for selection as coalescent model choice. *GENOME BIOLOGY AND EVOLUTION*, 12(6):871–877, JUN 2020.
- [35] Dennis Hedgecock and Alexander I. Pudovkin. Sweepstakes reproductive success in highly fecund marine fish and shellfish: a review and Commentary. *Bulletin of Marine Science*, 87(4):971–1002, OCT 2011.
- [36] Hussein A. Hejase, Ziyi Mo, Leonardo Campagna, and Adam Siepel. A deep-learning approach for inference of selective sweeps from the ancestral recombination graph. *MOLECULAR BIOLOGY AND EVOLUTION*, 39(1), JAN 7 2022.
- [37] Melissa Hubisz and Adam Siepel. Inference of ancestral recombination graphs using argweaver. In JY Dutheil, editor, *STATISTICAL POPULATION GENOMICS*, volume 2090 of *Methods in Molecular Biology*, pages 231–266. 2020.
- [38] RR HUDSON. Properties of a neutral allele model with intragenic recombination. *Theoretical Population Biology*, 23(2):183–201, 1983.
- [39] Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods. *Machine Learning*, 110(3):457–506, March 2021.
- [40] Ulas Isildak, Alessandro Stella, and Matteo Fumagalli. Distinguishing between recent balancing selection and incomplete sweep using deep neural networks. *Molecular Ecology Resources*, 21(8):2706–2718, November 2021.
- [41] Parul Johri, Charles F. Aquadro, Mark Beaumont, Brian Charlesworth, Laurent Excoffier, Adam Eyre-Walker, Peter D. Keightley, Michael Lynch, Gil McVean, Bret A. Payseur, Susanne P. Pfeifer, Wolfgang Stephan, and Jeffrey D. Jensen. Recommendations for improving statistical inference in population genomics. *PLOS Biology*, 20(5):e3001669, May 2022.
- [42] Parul Johri, Brian Charlesworth, and Jeffrey D. Jensen. Toward an evolutionarily appropriate null model: Jointly inferring demography and purifying selection. *GENETICS*, 215(1):173–192, MAY 2020.
- [43] Parul Johri, Kellen Riall, Hannes Becher, Laurent Excoffier, Brian Charlesworth, and Jeffrey D. Jensen. The impact of purifying and background selection on the inference of population history: Problems and prospects. *MOLECULAR BIOLOGY AND EVOLUTION*, 38(7):2986–3003, JUL 2021.
- [44] Mamoru Kato, Daniel A. Vasco, Ryuichi Sugino, Daichi Narushima, and Alexander Krasnitz. Sweepstake evolution revealed by population-genetic analysis of copy-number alterations in single genomes of breast cancer. *Royal Society of Open Science*, 4(9), SEP 2017.

- [45] Jerome Kelleher, Kevin R. Thornton, Jaime Ashander, and Peter L. Ralph. Efficient pedigree recording for fast population genetics simulation. 14(11):e1006581.
- [46] Jerome Kelleher, Yan Wong, Anthony W. Wohns, Chaimaa Fadil, Patrick K. Albers, and Gil McVean. Inferring whole-genome histories in large population datasets (vol 51, pg 1330, 2019). *Nature Genetics*, 51(11):1660, NOV 2019.
- [47] Younhun Kim, Frederic Koehler, Ankur Moitra, Elchanan Mossel, and Govind Ramnarayan. How Many Subpopulations Is Too Many? Exponential Lower Bounds for Inferring Population Histories. *Journal of Computational Biology*, 27(4):613–625, APR 1 2020.
- [48] JFC Kingman. The Coalescent . *Stochastic Processes and their Applications*, 13, 1982.
- [49] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks. 2016.
- [50] Kevin Korfmann, Oscar E Gaggiotti, and Matteo Fumagalli. Deep Learning in Population Genetics. *Genome Biology and Evolution*, 15(2):evad008, February 2023.
- [51] Jere Koskela. Multi-locus data distinguishes between population growth and multiple merger coalescents. *STATISTICAL APPLICATIONS IN GENETICS AND MOLECULAR BIOLOGY*, 17(3), JUN 2018.
- [52] Jere Koskela and Maite Wilke Berenguer. Robust model selection between population growth and multiple merger coalescents. *Mathematical Biosciences*, 311:1–12, MAY 2019.
- [53] John Boaz Lee, Ryan Rossi, and Xiangnan Kong. Graph Classification using Structural Attention. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1666–1674, London United Kingdom, July 2018. ACM.
- [54] Heng Li and Richard Durbin. Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357):493–U84, JUL 28 2011.
- [55] Ali Mahmoudi, Jere Koskela, Jerome Kelleher, Yao-ban Chan, and David Balding. Bayesian inference of ancestral recombination graphs. *PLOS COMPUTATIONAL BIOLOGY*, 18(3), MAR 2022.
- [56] P Marjoram and JD Wall. Fast “coalescent” simulation. *BMC Genetics*, 7, MAR 15 2006.
- [57] Sebastian Matuszewski, Marcel E. Hildebrandt, Guillaume Achaz, and Jeffrey D. Jensen. Coalescent processes with skewed offspring distributions and non-equilibrium demography. *Genetics*, 2017.
- [58] GAT McVean and NJ Cardin. Approximating the coalescent with recombination. *Philosophical Transactions of the Royal Society B-Biological Sciences*, 360(1459):1387–1393, JUL 29 2005.

- [59] F Menardo, S Gagneux, and F Freund. Multiple merger genealogies in outbreaks of *Mycobacterium tuberculosis*. *Molecular Biology and Evolution*, 07 2020. msaa179.
- [60] Alistair Miles, pyup io bot, Murillo R, Peter Ralph, Nick Harding, Rahul Pisupati, Summer Rae, and Tim Millar. cggh/scikit-allel: v1.3.3.
- [61] M Mohle and S Sagitov. A classification of coalescent processes for haploid exchangeable population models. *Annals of Probability*, 29(4):1547–1562, OCT 2001.
- [62] Ana Y. Morales-Arce, Rebecca B. Harris, Anne C. Stone, and Jeffrey D. Jensen. Evaluating the contributions of purifying selection and progeny-skew in dictating within-host *Mycobacterium tuberculosis* evolution. *Evolution*, 74(5):992–1001, MAY 2020.
- [63] Richard A. Neher and Oskar Hallatschek. Genealogies of rapidly adapting populations. *Proceedings of the National Academy of Sciences*, 110(2):437–442, January 2013.
- [64] Dominic Nelson, Jerome Kelleher, Aaron P. Ragsdale, Claudia Moreau, Gil McVean, and Simon Gravel. Accounting for long-range correlations in genome-wide simulations of large cohorts. *PLOS Genetics*, 16(5), MAY 2020.
- [65] Hiro-Sato Niwa, Kazuya Nashida, and Takashi Yanagimoto. Reproductive skew in japanese sardine inferred from dna sequences. *ICES Journal of Marine Science*, 73(9):2181–2189, 2016.
- [66] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. October 2017.
- [67] J Pitman. Coalescents with multiple collisions. *Annals of Probability*, 27(4):1870–1902, OCT 1999.
- [68] Xinghu Qin, Charleston W. K. Chiang, and Oscar E. Gaggiotti. Deciphering signatures of natural selection via deep learning. *bioRxiv*, 2021.
- [69] Matthew D. Rasmussen, Melissa J. Hubisz, Ilan Gronau, and Adam Siepel. Genome-wide inference of ancestral recombination graphs. *PLOS GENETICS*, 10(5), MAY 2014.
- [70] Daniel P Rice, John Novembre, and Michael M Desai. Distinguishing multiple-merger from kingman coalescence using two-site frequency spectra. *bioRxiv*, 2018.
- [71] Alan R. Rogers and Chad Huff. Linkage disequilibrium between loci with unknown phase. 182(3):839–844.
- [72] Andrew M. Sackman, Rebecca B. Harris, and Jeffrey D. Jensen. Inferring demography and selection in organisms characterized by skewed offspring distributions. *GENETICS*, 211(3):1019–1028, MAR 2019.
- [73] S Sagitov. The general coalescent with asynchronous mergers of ancestral lines. *Journal of Applied Probability*, 36(4):1116–1125, DEC 1999.

- [74] S Sagitov. Convergence to the coalescent with simultaneous multiple mergers. *Journal of Applied Probability*, 40(4):839–854, DEC 2003.
- [75] Théophile Sanchez, Jean Cury, Guillaume Charpiat, and Flora Jay. Deep learning for population size history inference: Design, comparison and combination with approximate bayesian computation. 21(8):2645–2660.
- [76] Nicolae Sapoval, Amirali Aghazadeh, Michael G. Nute, Dinler A. Antunes, Advait Balaji, Richard Baraniuk, C. J. Barberan, Ruth Dannenfelser, Chen Dun, Mohammadamin Edrisi, R. A. Leo Elworth, Bryce Kille, Anastasios Kyrillidis, Luay Nakhleh, Cameron R. Wolfe, Zhi Yan, Vicky Yao, and Todd J. Treangen. Current progress and open challenges for applying deep learning across the biosciences. *Nature Communications*, 13(1):1728, December 2022.
- [77] Ori Sargsyan and John Wakeley. A coalescent process with simultaneous multiple mergers for approximating the gene genealogies of many marine organisms. *Theoretical population biology*, 74(1):104–114, 2008.
- [78] Stephan Schiffels and Richard Durbin. Inferring human population size and separation history from multiple genome sequences. *Nature Genetics*, 46(8):919–925, AUG 2014.
- [79] Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks, 2017.
- [80] J Schweinsberg. Coalescent processes obtained from supercritical Galton-Watson processes. *Stochastic Processes and their Applications*, 106(1):107–139, JUL 2003.
- [81] Thibaut Paul Patrick Sellinger, Diala Abu Awad, Markus Moest, and Aurelien Tellier. Inference of past demography, dormancy and self-fertilization rates from whole genome sequence data. *PLOS Genetics*, 16(4), APR 2020.
- [82] Thibaut Paul Patrick Sellinger, Diala Abu-Awad, and Aurelien Tellier. Limits and convergence properties of the sequentially markovian coalescent. *MOLECULAR ECOLOGY RESOURCES*, 21(7):2231–2248, OCT 2021.
- [83] Sara Sheehan and Yun S. Song. Deep Learning for Population Genetic Inference. *PLOS Computational Biology*, 12(3), MAR 2016.
- [84] Leo Speidel, Marie Forest, Sinan Shi, and Simon R. Myers. A method for genome-wide genealogy estimation for thousands of samples. *Nature Genetics*, 51(9):1321+, SEP 2019.
- [85] Matthias Steinruecken, Matthias Birkner, and Jochen Blath. Analysis of DNA sequence variation within marine species using Beta-coalescents. *Theoretical Population Biology*, 87:15–24, AUG 2013.
- [86] Wolfgang Stephan. Selective Sweeps. *Genetics*, 211(1):5–13, January 2019.
- [87] Stefan Struett, Thibaut Sellinger, Sylvain Glémin, Aurélien Tellier, and Stefan Laurent. Inference of evolutionary transitions to self-fertilization using whole-genome sequences. *bioRxiv*, 2022.

- [88] Aurelien Tellier and Christophe Lemaire. Coalescence 2.0: a multiple branching of recent theoretical developments and their applications. *Molecular Ecology*, 23(11):2637–2652, JUN 2014.
- [89] Jonathan Terhorst, John A. Kamm, and Yun S. Song. Robust and scalable inference of population history from hundreds of unphased whole genomes. *Nature Genetics*, 49(2):303–309, FEB 2017.
- [90] Gautam Upadhyay and Matthias Steinrücken. Robust Inference of Population Size Histories from Genomic Sequencing Data. preprint, Genetics, May 2021.
- [91] Ke Wang, Iain Mathieson, Jared O’Connell, and Stephan Schiffels. Tracking human population structure through time from whole genome sequences. *PLOS Genetics*, 16(3), MAR 2020.
- [92] Zhanpeng Wang, Jiaping Wang, Michael Kourakos, Nhung Hoang, Hyong Hark Lee, Iain Mathieson, and Sara Mathieson. Automatic inference of demographic parameters using generative adversarial networks. *Molecular Ecology Resources*, 21(8):2689–2705, 2021.
- [93] C Wiuf and J Hein. Recombination as a point process along sequences. *Theoretical Population Biology*, 55(3):248–259, JUN 1999.
- [94] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [95] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. *CoRR*, abs/1603.08861, 2016.
- [96] Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec, Francesco Montinaro, Cyril Furtlehner, Luca Pagani, and Flora Jay. Creating artificial human genomes using generative neural networks. *PLOS Genetics*, 17(2):1–22, 02 2021.
- [97] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling.
- [98] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [99] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. 1:57–81.