# MOONSHINE.JL: A JULIA PACKAGE FOR GENOME-SCALE MODEL-BASED ANCESTRAL RECOMBINATION GRAPH INFERENCE

### A PREPRINT

**Patrick Fournier**
Département de Mathématiques
Université du Québec à Montréal
Montréal, Québec, H3C 3P8
p_fournier@hushmail.com

**Fabrice Larribe**
Département de Mathématiques
Université du Québec à Montréal
Montréal, Québec, H3C 3P8
larribe.fabrice@uqam.ca

November 27, 2025

### ABSTRACT

The ancestral recombination graph (ARG) is the model of choice in statistical genetics to model population ancestries. Software capable of simulating ARGs on a genome scale within a reasonable amount of time are now widely available for most practical use cases. While the inverse problem of inferring ancestries from a sample of haplotypes has seen major progress in the last decade, it does not enjoy the same level of advancement as its counterpart. Up until recently, even moderately sized samples could only be handled using heuristics. In recent years, the possibility of model-based inference for datasets closer to "real world" scenarios has become a reality, largely due to the development of threading-based samplers. This article introduces Moonshine, a Julia package that has the ability, among other things, to infer ARGs for samples of thousands of human haplotypes of sizes on the order of hundreds of megabases within a reasonable amount of time. On recent hardware, our package is able to infer an ARG for samples of densely haplotyped (over one marker/kilobase) human chromosomes of sizes up to 10000 in well under a day on data simulated by msprime. Scaling up simulation on a compute cluster is straightforward thanks to a strictly single-threaded implementation. While model-based, it does not resort to threading but rather places restrictions on probability distributions typically used in simulation software in order to enforce sample consistency. In addition to being efficient, a strong emphasis is placed on ease of use and integration into the biostatistical software ecosystem.

***Keywords*** Ancestral Recombination Graph · ARG Inference · Coalescent Theory · Algorithms · Julia

## 1 Introduction

Coalescent theory is the framework of choice when it comes to inference from genetic data ([1]). A brief introduction and some historical notes on the subject can be found in [2]. Although it originates from the seminal papers of Kingman ([3, 4]) and Hudson ([5]), it didn't really gain widespread use before the advent of computer software capable of simulating the coalescent process with recombination. ms ([5]), a software written in C, was the first of its kind. It aims at generating a sample of genetic sequences evolving under a Wright-Fisher (WF) model via a coalescent-with-recombination (CWR) approximation. The WF model is known for its simplicity. Somewhat surprisingly, simulating a sample of genetic sequences resulting from the evolution of a population obeying it rapidly becomes a challenging task. This is because it is necessary to keep track of every sequence at each of the numerous discrete steps of the simulation procedure from the original population to the desired sample. The coalescent process works in reverse: it starts from the sample and generates coalescence events, which correspond to sets of sequences finding a common ancestor. Consequently, it is only necessary to keep track of haplotypes associated with these events. The remaining sequences are deemed *non-ancestral* for the sample and can be disregarded. In addition, exponentially distributed

coalescence times are substituted to discrete generation counts. Sampling genetic sequences using the coalescent approximation is a straightforward task.

What is not trivial, however, is relaxing the assumptions of the WF model to allow recombination events. The result of simulations performed under a model accounting for these events, along with, somewhat confusingly, the statistical model itself, is known as the ancestral recombination graph (ARG). To avoid ambiguity, we reserve "ARG" for the former and refer to the statistical model through "CWR". While Hudson's ms is able to simulate recombination events, it is not efficient enough to meet the requirements of large-scale genomic data analysis. To reduce the computational burden associated with simulations, an approximation to the CWR was developed by Niall Cardin and Gil McVean ([6]). Building on the work of [7], their idea is to approximate the non-Markovian recombination process with a Markovian one called the sequential Markov coalescent (SMC). This gave rise to a plethora of so-called sequential simulators such as MaCS ([8]), fastsimcoal ([9]), SC ([10]) and scrm ([11]). While the sequential approach is not inherently more efficient than the backward-in-time one, it is easier to approximate, leading to substantial improvements in the capacity of CWR samplers. Nonetheless, the idea of backward-in-time simulation was not abandoned, as evidenced by MSMS ([12]) and Discoal ([13]). [14] even developed a forward-in-time simulator able to handle non-WF scenarios. Some years later, the classical backward-in-time CWR approach was brought back into the spotlight by msprime ([15]). By reformulating Hudson's original ms in terms of a data structure they called *sparse trees*, their Python package (with performance-critical procedures implemented in C) has the capacity of simulating exactly from the coalescent with recombination more efficiently than the sequential approximations available at the time. Version 1.0 ([16]) brings even more features such as the ability to simulate directly from a WF model. For all these reasons, msprime is nowadays the *de facto* reference for CWR and even ARG simulation.

Programs mentioned above simulate genealogies with the goal of generating a sample of genetic sequences. Throughout this paper, we refer to this process as *ARG simulation*; the ancestry is viewed as a sample point in the probability space associated with the CWR. The focus is generally on inferring parameters such as recombination/mutation rates or effective population size. The distinctive characteristic of this class of samplers is that the distribution is parametrized by values derived from sequences of markers rather than the markers themselves. This is to be contrasted with software designed to solve the inverse problem of generating *likely* ancestries, or even, in some cases, a single ancestry directly from the markers. This is known as *ARG inference* since the ancestry itself is usually the main point of interest. This choice of words is, however, somewhat misleading as so-called inferred genealogies are not necessarily central to the analyses they are involved in. They can be instrumental in the estimation of other parameters, just like their sampled counterparts. It is true that their distribution is not that of the CWR. That being said, this is not to say that alternative likelihood functions enabling maximum likelihood estimation do not exist. Indeed, recent work ([17]) proposes formulations with applicability to broad classes of ancestries in mind. While the problem these samplers solve is by nature more computationally demanding, they provide major benefits in that the ancestries they produce are generally more likely with respect to the sample at hand. In particular, their usefulness cannot be overstated for methodologies that aim at improving the estimation of parameters by treating the ancestry of a sample as a latent variable such as [18] or [19]. Those require integrating these genealogies out, a task involving, in practice, the ability to sample in high probability regions.

Early attempts at ARG inference include recom ([20]) and Infs ([21]). More recently, the authors of SC also provides a modified version of their algorithm called *SC-sample*, capable of solving this inverse problem by generating sample-consistent graphs. A graph is said to be consistent for a sample with respect to a mutation evolution model if it generates the sample under that model. In [10], they choose to use the popular infinite site model (ISM). Other software have been developed with the goal of inferring the ancestry of a sample, such as tsinfer ([22]) and ARGinfer ([23]). Both assume a sample of sequences of binary markers evolving under the ISM. tsinfer uses a heuristic to partially infer haplotypes of unobserved ancestors. Each site of every sequence in the sample is then traced back to an ancestor by maximum likelihood. This yields a point estimate of the ancestral recombination graph. The goal of ARGinfer is slightly different, as it aims at performing the inference probabilistically. It can compute probability intervals for ARG related quantities. All these simulators are designed to sample ARGs according to a probabilistic model. For that reason, we refer to them as *model-based simulators*. An extensive review of ARG samplers is available in [24].

A related problem is that of *parsimony-based inference*, which consists of finding ARGs consistent with a sample using the minimum number of recombination events. This problem has been proved NP-hard ([25]). Nonetheless, attempts to solve it exactly and approximately go back to the early work of Hein ([26]), which was improved by [27]. Other algorithms include Margarita ([28]), ARG4WG([29]) and GAMARG([30]).

This paper introduces a Julia package called Moonshine implementing inference of sample-consistent ancestral recombination graphs. The approach is sequential, as it is based on iterative modification of the ARG. A sequence of operations are applied to a coalescent tree to ultimately make it consistent with a sample of haplotypes. The main advantage over back-in-time approaches is the possibility it offers users to specify various levels of approximation

when generating ARGs. The available spectrum ranges from *exact* simulation without any Markov assumption to first-order approximation *à la* SMC. The end result of the sampling routine is independent of approximation level. The whole graph, as well as meta-data such as vertex latitudes, associated haplotypes, and intervals of ancestrality for edges, is available to users for subsequent analysis. In fact, `Moonshine` is designed for easy integration into data analysis workflows. In addition, it is possible to use it solely for inferring a set of ancestries consistent with a given sample of genetic sequences. From that point of view, the software it is most similar to is probably `SC-sample`. That being said, both packages differ in many respects. At first sight, the most obvious difference is that `SC-sample` is a `C++` standalone software while `Moonshine` is a `Julia` package. The latter is much more flexible. A run of `SC-sample` results in a text file containing sequences of Newick-formated trees representing an ARG which may or may not be easily integrated into an existing pipeline. In contrast, thanks to `Julia`'s just-in-time compilation, `Moonshine` can be used interactively without sacrificing performance. It is also fully integrated with `Julia`'s ecosystem of graph theoretical packages. Another noteworthy difference is at the algorithmic level. `SC-sample` is a modified version of `SC` which is itself a modified version of `MaCS`. Although all sequential algorithms are based on the same idea and, consequently, share many similarities, `Moonshine` implements its own original algorithm. Furthermore, since it is created with statistical inference in mind, `Moonshine` treats ARGs as random graphs. It is straightforward to evaluate ARG-related functions such as probability densities for the ARGs themselves or other random variables, such as phenotypes, conditional on an ARG. Implementation of custom functionalities is facilitated by a coherent type hierarchy and thorough documentation of abstract types and interfaces, making the extension of the package's various components as easy as possible. Finally, interoperability with `tskit` ([15, 31]) streamlines data management. Integration with the popular genomics package is still a work in progress, but creating samples directly from `tskit` objects is supported. A convenience constructor for generating random samples from a simple genetic model using `msprime` ([16]) is provided. This is transparent to the end user, thanks to `Moonshine` being packaged with its own distribution of `msprime`.

Our objective in developing `Moonshine`'s is not limited to creating a realistic and convenient sampler; performance is a major priority. We present numerical experiments showing its potential for both coalescent tree construction and ARG inference. Trees for large samples ($n = 10000$) of long simulated haplotypes (250 Mbp) can be constructed in minutes at high resolution (over one marker per kbp) using Hamming distance between sequences. In the same scenarios, complete ancestries can be inferred in hours. Furthermore, our algorithms are completely single-threaded, enabling us to increase sampling throughput by leveraging concurrency efficiently and easily by launching multiple instances in parallel.

## 2 Tree Construction

Similar to other sequential samplers, the first step of our algorithm is to construct an initial coalescent tree. Consistency with the first marker is not assumed; ARG inference can be carried out even starting from a completely random tree. The sole requirement is that it be a valid coalescent tree, i.e., a full binary tree with a coherent set of latitudes for the vertices. The idea behind this functionality is that it might be of interest to compare the performance of a method under a model without recombination versus one that allows for such events. It would make little sense in that context to give a special status to a single marker, disregarding the remainder of the haplotypes. Consequently, we give the user maximum flexibility when building coalescent trees, which may be of interest both in their own right and as a stepping stone for constructing more complex histories. As will be discussed later, `Moonshine` is packaged with two haplotype metrics designed for tree building. It is straightforward for the user to implement custom metrics.

Within our package, data structures representing genealogies are subtypes of the `AbstractGenealogy` abstract type. The type of coalescent trees is simply `Tree`. Given a sample (of type `Sample`) of phased and polarized genetic sequences (of type `Sequence`), sampling is controlled by two parameters: the global mutation rate $\mu$ and a metric on haplotypes $d$. Assuming a diploid population, sequences of length $l$ with $s$ markers, an effective population size $N_e$ and a constant per locus mutation rate of $\mu'$, the global mutation rate is $\mu = 2N_e\mu'l$. These parameters are either computed or explicitly passed to `Sample`'s constructor. As for the metric, since `Moonshine` is compatible with binary markers exclusively, a sample is an $n$-tuple $\boldsymbol{H} = (h_1, \cdots, h_n)$ of $s$-vectors of GF(2), the finite field with two elements. Concretely, for each $k$, we have

$$h_k = h_k^1 \cdots h_k^s$$

where $h_k^\bullet \in \{0, 1\}$. Wild and derived alleles are represented by 0 and 1 respectively as is standard in the literature. Let $\oplus$ denote addition modulo 2. Examples of useful metrics, also known as *distance functions*, include

$$d_L(h_1, h_2) = h_1^1 \oplus h_2^1$$

$$d_H(h_1, h_2) = \sum_{i=1}^s h_1^i \oplus h_2^i \,.$$

3

$d_L$ is the metric under which the distance between two sequences is zero if and only if the state of their first marker is identical, 1 otherwise. $d_H$ is the Hamming distance. Such discrete distances have a natural biological interpretation as the number of mutations between (a subinterval of) sequences. Arbitrary distances can be implemented by the user as subtypes of `Distance`.

Detail-oriented readers might have noticed that the term "metric" is used loosely, as the positivity axiom need not hold; the distance between two distinct haplotypes may be 0. This is necessary to allow inference of trees consistent for a single marker. Technically, the correct mathematical construct is that of a pseudometric.

## 2.1 Sampling Algorithm

Coalescence events are sampled as follows: a vertex $v_a$ is choosen uniformly among the set of *live* vertices, that is, the vertices that have not coalesced yet. Another vertex $v_b$ is sampled conditional on $v_a$. The probability $p_{ab}$ of sampling $v_b$ given $v_a$ is proportional to

$$p_{ab} = \frac{\mu^{d_{ab}}}{\Gamma(d_{ab} + 1)} \tag{1}$$

where $\Gamma$ is the gamma function, $d_{ab} = d(h_a, h_b)$ with $h_a$ and $h_b$ the haplotypes associated with $v_a$ and $v_b$ respectively. $p_{ab}$ is an unnormalized Poisson probability, where $\Gamma$ is used instead of the usual factorial function to allow non-integer distances. Both vertices then coalesce into $v_c$ with associated haplotype $h_c = h_a \odot h_b$ where $\odot$ is the Hadamard product. The shift in latitude is exponentially distributed with rate parameter equal to the number of live vertices. The latitude of $v_c$ is computed with respect to that of the previous event. This procedure replaces $v_a$ and $v_b$ by $v_c$ in the set of live vertices. Repeating it $n-1$ times on a sample of $n$ haplotypes yields a coalescent tree.

The user has the possibility of biasing distance computation by specifying a parameter $c_0 \in \mathbb{R}_+ \cup \{0, \infty\}$ meaning that, in practice,

$$d_{ab} = c_0 d(h_a, h_b).$$

$c_0$ can be used to reproduce the behavior of other samplers. Setting $c_0 = \infty$ and $d = d_L$, we obtain

$$d_{ab} = \begin{cases} 0 & \text{if } h_a^0 = h_b^0 \\ \infty & \text{if } h_a^0 \neq h_b^0 \end{cases}$$

resulting in

$$p_{ab} = \begin{cases} 1 & \text{if } h_a^0 = h_b^0 \\ 0 & \text{if } h_a^0 \neq h_b^0 \end{cases}.$$

The resulting sampler will aggregate haplotypes with identical status at the first marker, resulting in a tree consistent with it.

In practice, dealing with a ratio of such extreme quantities poses a numerical challenge. For instance, if one wishes to use Hamming's distance, computation of the normalizing constant becomes impossible even for relatively small values of $n$ and $s$. We address these issues using two tricks. First, we compute $p_{ab}$ via Stirling's approximation. On the logarithmic scale, we obtain

$$\log p_{ab} \approx d_{ab}(\log \mu - \log d_{ab} + 1)$$

which is already more manageable. Next, it would be ideal to sample without reverting to the linear scale. Moreover, we would greatly benefit from avoiding the computation of the normalizing constant altogether. It turns out that this is exactly what the so-called *Gumbel trick* [32] is designed to do. The trick transforms sampling from the target distribution into an optimization problem. For a set of candidate vertices $b_1, \ldots, b_m$ and a sequence of iid standard Gumbel random variables $t_1, \ldots, t_m$, $\arg\max_k\{\log p_{ab_k} + t_k\}$ is distributed as a categorical random variable with the probability of category $k$ being equal to $p_{ab_k}$. This result gives us a more stable way of sampling the second coalescing vertex $v_b$, at the cost of increased computing time for drawing Gumbel random variables and finding the maximum of the sequence. In many cases, this impact should be minimal compared to the overall execution time, for instance when the tree is to be used for ARG inference. Figure 1 provides a rough estimate of time and memory usage for various scenarios. For cases where time considerations warrant a precision tradeoff, users can sample approximately from the target distribution via a scheme we call *secretary sampling*, inspired by the famed *secretary problem* ([33]). It revolves around giving the algorithm a chance of terminating before having traversed the complete set of candidate vertices. Its behavior is controlled by a user-determined threshold parameter $t_0 \in [0, 1]$. As illustrated in fig. 2, the probability of an early termination decreases with $t_0$. Note that although normalization is not required for sampling, the unnormalized probabilities associated with traversed vertices must be summed in order to evaluate the density of the resulting coalescent tree. This has to be done carefully as departure from the logarithmic scale is unavoidable. Details as well as the complete sampling algorithm are presented in algorithm 1.

---

**Algorithm 1:** Tree Sampling

---

**1**  Fill array 'live' with the live edges

**2**  **while** $|\mathsf{live}| > 1$ **do**

**3**       Sample $\mathsf{v_a}$ uniformly from $\mathsf{live}$ and remove it from $\mathsf{live}$

        // $\log p_{ab} + \mathtt{Gumbel}$, $\log p_{ab}$, `normalization constant and index of sampled vertex`
        `respectively`

**4**       Set $\mathsf{z} \leftarrow \infty$, $\mathsf{logp} \leftarrow 0$, $\mathsf{logc} \leftarrow 0$ and $\mathsf{idx} \leftarrow 0$

**5**       Set $\mathsf{k} \leftarrow 0$

**6**       **while** $\mathsf{k} < |\mathsf{live}|$ **do**

**7**           Increment $\mathsf{k} \leftarrow \mathsf{k} + 1$

**8**           Set $\mathsf{b} \leftarrow \mathsf{live[k]}$, $\mathsf{logp} \leftarrow \log p_{ab}$

**9**           **if** $\mathsf{logp} < \infty$ **then**

**10**              Sample $\mathsf{g} \sim \mathtt{Gumbel}(\textit{0, 1})$

**11**              Set $\mathsf{logc} \leftarrow \mathsf{logp}$, $\mathsf{z} \leftarrow \mathsf{logc} + \mathsf{g}$ and $\mathsf{idx} \leftarrow \mathsf{k}$

**12**              Break

**13**       **while** $\mathsf{k} < |\mathsf{live}|$ **do**

**14**           Increment $\mathsf{k} \leftarrow \mathsf{k} + 1$

**15**           Set $\mathsf{b} \leftarrow \mathsf{live[k]}$, $\mathsf{logp'} \leftarrow \log p_{ab}$

**16**           **if** $\mathsf{logp'} < \infty$ **then**

            // `Update` $\mathsf{logc}$ `accurately`

**17**              Set $\mathsf{pmin} \leftarrow \min\{\mathsf{logp'}, \mathsf{logc}\}$, $\mathsf{pmax} \leftarrow \max\{\mathsf{logp'}, \mathsf{logc}\}$

**18**              Update $\mathsf{logc} \leftarrow \mathsf{pmax} + \log(1 + \exp(\mathsf{pmin} - \mathsf{pmax}))$

**19**              Sample $\mathsf{g} \sim \mathtt{Gumbel}(\textit{0, 1})$

**20**              Set $\mathsf{z'} \leftarrow \mathsf{logp'} + \mathsf{g}$

**21**              **if** $\mathsf{z'} \geq \mathsf{z}$ **then**

**22**                  Set $\mathsf{logp} \leftarrow \mathsf{logp'}$, $\mathsf{z} \leftarrow \mathsf{z'}$, $\mathsf{idx} \leftarrow \mathsf{k}$

**23**                  **if** $\mathsf{k} > \textit{threshold}$ **then**

**24**                     Break

**25**       **if** $\mathsf{idx} = 0$ **then** // `All probabilities were infinite`

**26**          Sample $\mathsf{idx}$ uniformly from $1, \dots, |\mathsf{live}|$

**27**          Set $\mathsf{logp} \leftarrow 0$, $\mathsf{logc} \leftarrow \log |\mathsf{live}|$

**28**       Set $\mathsf{b} \leftarrow \mathsf{live[idx]}$, $\mathsf{v_c} \leftarrow 2n - |\mathsf{live}|$

**29**       Remove $\mathsf{v_b}$ from $\mathsf{live}$, coalesce $\mathsf{v_a}$ and $\mathsf{v_b}$ into $\mathsf{v_c}$ and add $\mathsf{v_c}$ to $\mathsf{live}$

**30**       Add $\mathsf{logp} - \mathsf{logc} - \log(1 + |\mathsf{live}|)$ to tree's log density
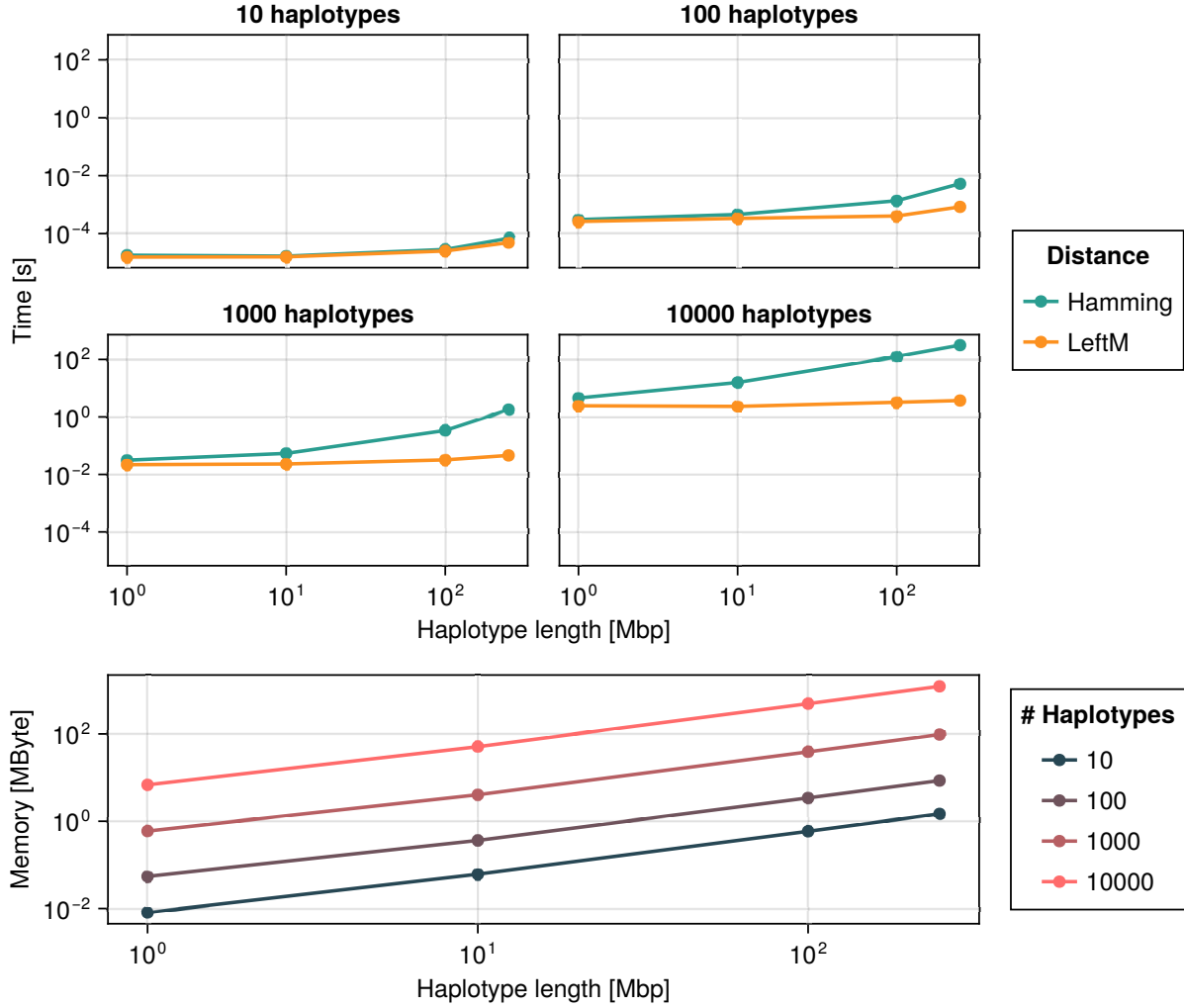
---

Figure 1: Time and memory needed for the construction of a coalescent tree as a function of the number of haplotypes, haplotype length, and distance function. Sampling is exact ($t_0 = 1$) and no bias is applied ($c_0 = 1$). Sampling with respect to either of the two distances yields identical memory usage since their computation does not involve memory allocation. Results are also available in table 1.

## 3  ARG Inference

In `Moonshine`, ancestral recombination graphs are instances of the `ARG` type, a subtype of `AbstractGenealogy`. Since they represent a more realistic model for the ancestry of a sample subject to recombination, ARGs can be viewed as improved versions of coalescent trees. `ARGs` are constructed sequentially from a `Tree` by iteratively sampling recombination events until an ancestry consistent with the sample is reached. As is common, consistency is defined through the ISM: a genealogy is consistent with a sample if the number of mutations per marker is at most one. In a context where the recombination rate is several orders of magnitude higher than the mutation rate, a consistent ARG is typically a more realistic genealogy than any other kind of inconsistent ancestry. `Moonshine` is consequently very well suited to working with single nucleotide polymorphisms (SNPs), which have a low mutation rate.

Ancestries are modified by recombination events, which partitions ancestral material into two subintervals: that to the left of an associated point, called a *breakpoint*, and the material to the right. From a graph-theoretical perspective, a recombination event is generally represented by a vertex of degree 3 with a single child and two parents. As an example, imagine that the child edge of a recombination vertex is ancestral for an interval $I$ and that the associated breakpoint is $b$. In that case, one of the parental edges, generally referred to as the left edge, is ancestral for $[0, b) \cap I$
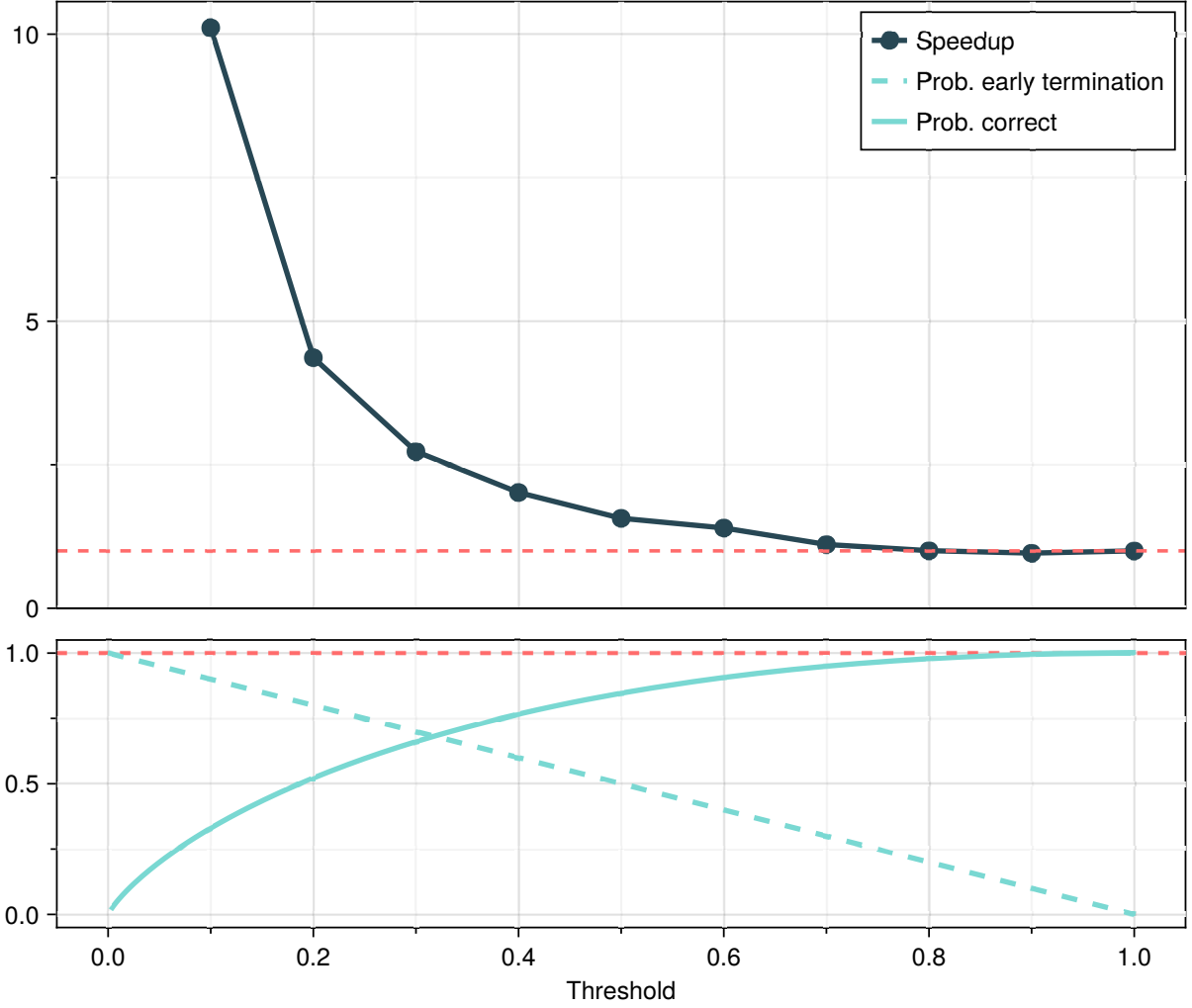
Figure 2: Speedup in tree construction, probability of sampling the correct sequence and probability of early termination as a function of the sampling threshold $t_0$. No bias is applied ($c_0 = 1$). The probability of early termination is $1 - t_0$. As the number of candidate vertices increases, the probability of sampling correctly from the target distribution converges to $t_0(1 - \log t_0)$ (depicted here). The exact probability for a single run is given by lemma 1. Results are also available in table 2.

while the other (right) edge is ancestral for $[b, \infty) \cap I$. Although we used the term "interval", $I$ can actually be a union of intervals. We will continue to use this terminology when the distinction between the two concepts is not relevant.

Recombination vertices are "added", figuratively speaking, by deleting an edge from the graph and connecting the new vertex with both endpoints of the removed edge. The edge connected to the child vertex is the recombination vertex's child edge, and the other is its left edge. This procedure leaves the right edge floating. Every recombination is immediately followed by the coalescence of the right edge with the graph. As it is carried out by a different algorithm than the coalescence of two dangling vertices encountered in temporal algorithms or, more trivially, when inferring a coalescent tree, we call those *recoalescence* events even though they result in an additional coalescence vertex as well. Coalescence and recombination vertices are in many regards mirror images of each other. A coalescence vertex has two child edges and one parental edge. If the child edges are ancestral for two intervals $I_l$ and $I_r$, then so is the parental edge for $I_l \cup I_r$. The standard recoalescence procedure begins again by deleting an edge, followed by connecting its incident vertices with the recoalescence vertex. The remaining child edge is then connected to the recombination vertex's right edge, concluding the procedure. Recoalescence can also occur without edge deletion. In that case, the coalescence vertex becomes the new root of the graph and lacks a parental edge. It is connected downstream to the previous root and the recombination vertex.

The type of a recoalescence event depends on its latitude, denoted $l_c$, which itself depends on the recombination's latitude $l_r$. The root of an ARG corresponds to the sample's most recent common ancestor (MRCA), sometimes called the *grand* MRCA (gMRCA). Its latitude is the time to the gMRCA (TgMRCA). Let $v_r$ and $v_c$ be the recombination and recoalescence vertices and $s_r - d_r$ and $s_c - d_c$ the edges deleted in the recombination-and-recoalescence (RR) steps. Algorithm 2 summarizes the RR procedure.
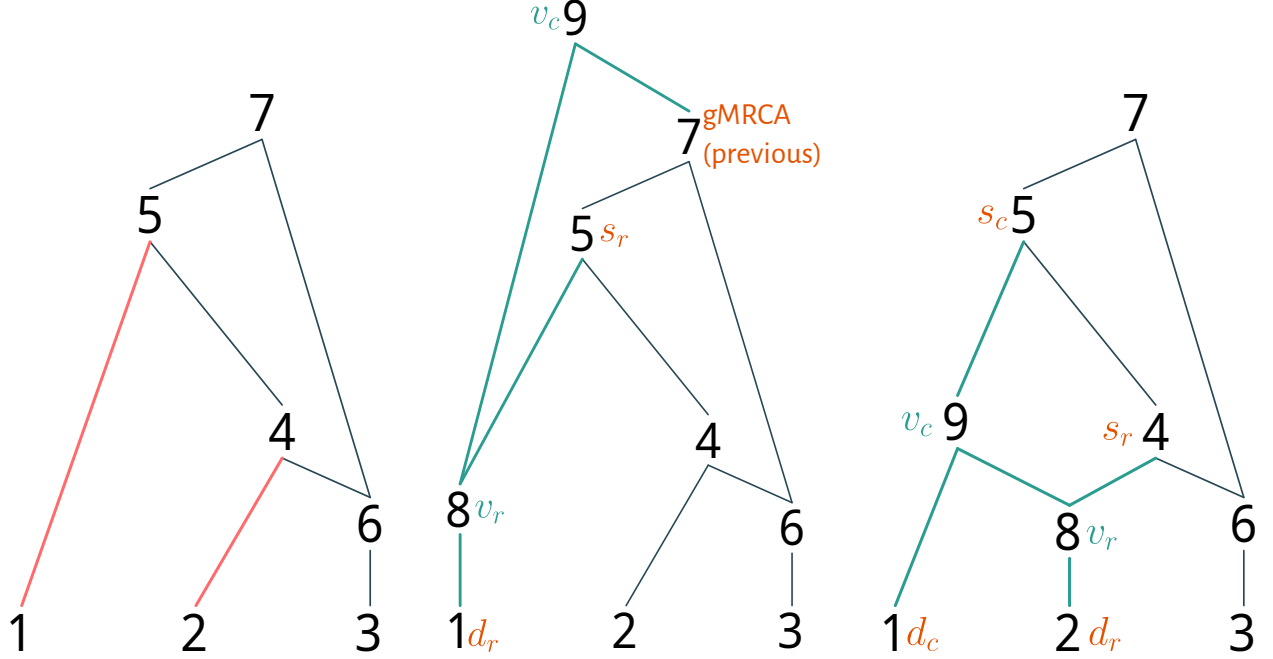
---

**Algorithm 2:** Recombination-and-recoalescence

**1** Add two vertices $\mathsf{v_r}$ and $\mathsf{v_c}$ to $\mathsf{G}$
**2** Delete recombination edge $\mathsf{s_r}$ — $\mathsf{d_r}$
**3** Add edges $\mathsf{s_r}$ — $\mathsf{v_r}$ and $\mathsf{v_r}$ — $\mathsf{d_r}$
**4** **if** $\mathsf{l_c} \leq \mathsf{TgMRCA}$ **then**
**5**      Delete recoalescence edge $\mathsf{s_c}$ — $\mathsf{d_c}$
**6**      Add edges $\mathsf{s_c}$ — $\mathsf{v_c}$ and $\mathsf{v_c}$ — $\mathsf{d_c}$
**7** **else**
**8**      Add edge $\mathsf{v_c}$ — $\mathsf{mrca}$
**9** Add edge $\mathsf{v_c}$ — $\mathsf{v_r}$

---

### 3.1 Unrestricted Recombination-and-Recoalescence events

`Moonshine` has the capability to sample two kinds of RR events: restricted and unrestricted. Unrestricted events, the subject of this section, have a distribution designed to closely match that of the CWR. Restricted recombination events are sampled to reduce the total number of mutation events on an ARG; these will be discussed at length in the next section.

Our package exports a method for sampling an arbitrary number of unrestricted recombination events. Standard theory [7] models the positions (on sequences) and locations (on ancestries) of recombination events as a Poisson point process (PPP). Conditional on their number, both locations and positions are distributed uniformly. When applied directly, this method has the drawback of generating sequences devoided of material ancestral for the sample at hand. Recombination events can be classified depending on whether ancestral material is present on both sides. If it is, the breakpoint can be positioned in ancestral or non-ancestral material. These are referred to as *type 1* and *type 2* events respectively and exclusively create haplotypes having ancestral material. Events positioned such that only material to their left or right side is ancestral are classified as *type 3* and *type 4* respectively. Finally, events occurring in entirely non-ancestral sequences are classified as *type 5*. It is desirable for an algorithm to only generate the first two types of recombination events since the other ones do not contribute to the structure of the sample. Our method follows this approach. We start by drawing a recombination edge with probability proportional to its length. The location of the event is distributed conditional on the selected edge. Then, we sample a position, also known as a *breakpoint*, uniformly on the mathematical closure of the set of intervals for which the recombination edge is ancestral. This strategy enforces the existence of ancestral material on both sides of the breakpoint and avoids recombination events of type 3, 4 and 5. Additionally, since the closure of ancestral intervals is not, in general, equal to the intervals themselves (it may contain "holes" of non-ancestral material), type 2 events are possible.

(a) Initial graph. $4 - 2$ is the recombination edge in fig. 3c. $5 - 1$ is the recombination edge in fig. 3b and recoalescence edge in fig. 3c.

(b) Recoalescence above the TgMRCA. The resulting ARG is taller than the original.

(c) Recoalescence below the TgMRCA.

Figure 3: Two types of recombination events. Elements removed from the initial graph (fig. 3a) are in red while those added to it are in green.

Selecting a conditional distribution for the location of recombination events on branches requires careful consideration. Theory dictates that it should be uniform. However, this results in a problematic frequency of locations close to the branches' endpoints, contributing to the *short branches issue* discussed in subsection 3.2.3. Consequently, we use the location-scale family associated with the $\text{Beta}(2, 2)$ distribution instead. This method preserves the symmetry of the uniform distribution while favoring locations closer to the center of branches, achieving sufficient reduction in the number of short branches to prevent numerical error.

The recoalescence process follows standard theory. The latitude is distributed as an inhomogeneous Poisson process with rate equal to the number of branches. The usual time-scale transformation strategy (see subsection 3.2.3) is implemented. A recoalescence edge is sampled uniformly, conditional on the latitude. The process of sampling a single unrestricted recombination event is summarized in algorithm 3.

---
**Algorithm 3:** Unrestricted Recombination

---
1 Sample a recombination edge in $\mathsf{e_r} \in E$ with probability proportional to its length
2 Sample a breakpoint $\mathsf{r} \in (r_L, r_U)$ where $r_L$ and $r_U$ are the leftmost and rightmost positions for which $\mathsf{e_r}$ is ancestral
3 Sample a recombination latitude $\mathsf{l_r}$ distributed as $\mathtt{ldst}(\mathsf{e_r}) + (\mathtt{lsrc}(\mathsf{e_r}) - \mathtt{ldst}(\mathsf{e_r})) \, \text{Beta}(2, 2)$ where $\mathtt{lsrc}(\mathsf{e_r})$ and $\mathtt{ldst}(\mathsf{e_r})$ are the latitudes of the highest and lowest vertex adjacent to $\mathsf{e_r}$ respectively
4 Sample a recoalescence latitude $\mathsf{l_c}$ distributed as an inhomogeneous PPP with rate equal to the number of branches via time-scale transformation
5 Sample a recoalescence edge $\mathsf{e_c}$ uniformly among the edges at latitude $\mathsf{l_c}$
6 Apply algorithm 2

---

As a sidenote, since a function allowing users to add arbitrary recombination events to a graph is exposed, `Moonshine` can easily be used to simulate ancestries. We do not recommend doing so, however, especially not using the built-in types `Tree` and `Arg`, which are designed to implement ARG inference. Instances of these types store information irrelevant in a simulation context which greatly hinders performance.

## 3.2 Restricted Recombination-and-Recoalescence Events

`Moonshine` can be used as an inferred ARG sampler due to its ability to transform an inconsistent graph into a consistent one by generating a sequence of RR events. A typical run starts with sampling a coalescent tree and is followed by a sequential sweep made up of the following broad steps:

1. Start at the leftmost position;

2. Find the next position that is inconsistent with the sample;

3. Sample a restricted recombination event;

4. Repeat until the rightmost position is reached.

To ensure our graphs' distribution closely resembles that of the CWR, our aim, alongside computational efficiency, is to impose as few constraints as possible on restricted recombination events. The probability distributions involved are similar to their unrestricted counterparts, except for their support, which we attempt to reduce only as much as is necessary to enforce consistency of the final ARG.

### 3.2.1 SIMD-Accelerated Minimum Mutation Number Algorithm

Step 2 requires an efficient inconsistency detection algorithm. Specifically, given a starting position, it must be able to find the closest marker to its right that mutates more than once in its current marginal tree. Furthermore, efficiency considerations demand that the set of edges on which these mutation events occur, which we refer to as *mutation edges*, is computed simultaneously. We refer to algorithms fulfilling these two requirements as *minimum mutation number* (MMN) algorithms. Mutation edge identification is at the core of Moonshine's ARG inference procedure and it should come at no surprise that a lot of time and effort went into its optimization. It is key to enabling inference on real-world sized datasets. We begin this section by giving a general description of the procedure before delving into implementation specifics.

Our algorithm leverage the parallel introduced in section 2 between sequences of $s$ biallelic markers and $s$-vectors of $\mathrm{GF}(2)$. In addition to $\oplus$ defined earlier, denote the multiplication on $\mathrm{GF}(2)$, which is identical to the one for natural numbers, by $\otimes$. Within that framework, both binary operations can be given a biological interpretation. Let $h_1$ and $h_2$ be two vectors of $\mathrm{GF}(2)^s$ and let $h = h_1 \oplus h_2$ where $\oplus$ is applied elementwise. Since the characteristic of $\mathrm{GF}(2)$ is 2 (i.e. $1 \oplus 1 = 0$), any non-zero element $h^{(0)}$ of $h$ results from $h_1^{(0)} \neq h_2^{(0)}$. This, in turn, indicates an odd number of mutations occurring between the two sequences at that marker. Assuming a mutation model disallowing back mutations, this simplifies to exactly one mutation event.

Interpretation of $\otimes$ is slightly more convoluted. Let $v$ be an internal vertex and $v_1, \ldots, v_m$ the subset of leaves having $v$ as an ancestor. Denote the corresponding sequences by $h_v, h_{v_1}, \ldots, h_{v_m}$. Assuming no back mutations,

$$h_v = h_{v_1} \otimes \cdots \otimes h_{v_m}$$

where, again, $\otimes$ is applied elementwise. This is easily seen to be the case for a vertex $v$ having two leaves $v_1$ and $v_2$ as children. The only way $h_v^k = 1$ is if $h_{v_1}^k = h_{v_2}^k = 1$ since at most one $k$-mutation could have occurred on $v - v_1$ and $v - v_2$. For the same reason, if $h_v^k = 0$, then either $h_{v_1}^k = h_{v_2}^k = 0$, in which case no mutation occurred, or $h_{v_1}^k = h_{v_2}^k \oplus 1$, in which case a single one occurred. This explanation can be extended to "deeper" vertices by moving upward, labelling internal vertices encountered along the way.

The correspondence between $\mathrm{GF}(2)$ and the evolution of sequences of biallelic markers under a model precluding back mutations can be exploited for computational efficiency. In practice, for large samples, our implementation reduces the time needed to compute mutation edges to about 15% of total execution on recent hardware, and this is without resorting to any form of multithreading. Nearly all of the time spent in the MMN algorithm is devoted to graph traversal. This level of performance is made possible by the following data structure: A sequence $h = h^1 \cdots h^s$ can be associated to the number $\widetilde{h}$ whose binary representation has its first $s$ least significant digits corresponding to $h$. In order to make this representation unique, we assume that the other digits of $\widetilde{h}$ are 0. This corresponds to

$$\widetilde{h} = \sum_{k=1}^{s} 2^k h^k \, .$$

When context is unambiguous, we shall drop the tilde and denote by $h$ either a haplotype, the corresponding vector or its integer representation. Treating haplotypes as integers is very convenient as operations $\oplus$ and $\otimes$ described above correspond to the bitwise exclusive disjunction (`XOR`) and conjunction (`AND`). For that reason, we denote those by $\oplus$

and $\otimes$ as well. This interpretation allows for very efficient handling of sequences since these are directly executed by the processor: we cannot break down mutation finding any further. Given an edge $v_1 - v_2$ ancestral for an interval $\omega$, we can find all of its mutations via algorithm 4:

---

**Algorithm 4:** Mutation Events on an Edge

---

**1** Compute a $\otimes$-mask $m_\omega$ masking markers outside $\omega$ to 0
**2** Compute $h = (h_{v1} \oplus h_{v_2}) \otimes m_\omega$
**3** Find the indices of set bits in $h$

---

Modern computers generally read data from memory in chunks of 8 bytes and can apply elementary bitwise operations such as `XOR` and `AND` to such chunks in a single cycle. This is readily implemented and results in the concurrent assessment of 64 markers on two haplotypes on a single physical core. This is good, but we can do even better. On recent hardware, some operations can be applied to chunks larger than 8 bytes, a capability enabled by so-called single instruction, multiple data (SIMD) instructions. At the upper end of the spectrum are processors supporting the AVX-512 instruction set, which can operate on 64 bytes simultaneously. Such operations are said to be "vectorized" in reference to them being applied to multiple 8-byte "scalars" at once. While not all regular CPU operations have a vectorized counterpart, integer `AND` and `XOR`, the only two operations required for mutation detection, are part of the AVX-512F extension available on any AVX-512 compliant processor.

That being said, it would generally not be very efficient to compare chunks of 512 markers from 2 different haplotypes simultaneously, as the next mutating marker is likely near the last recombination breakpoint. Instead, we process haplotypes in chunks of 8 markers. On an AVX-512 enabled architecture, this results in up to 64 haplotypes being compared at 8 positions concurrently without multithreading. In fact, this is so efficient that tested multithreaded versions were actually slower than the single-threaded one. Our algorithm is implemented in `Moonshine` using `SIMD.jl` ([34]) which ensures portability across architectures.

Computing the number of mutations by batches of 8 markers naturally leads to optimizations of the graph traversal procedure. The main one is based on the observation that, assuming no back mutation, a given marker cannot have the derived allele for a haplotype if one of its descendants has the wild allele. This is the biological version of 0 being the absorbing element of $\otimes$. By performing traversal in a bottom-up fashion starting at the leaves, we allow for early termination when encountering a vertex associated with a chunk of zeros. The probability of early termination increases as the size of the chunk of markers decreases, further justifying our decision to work with byte-sized chunks. Besides, bottom-up traversal is slightly more efficient than top-down traversal because the validity check associated with returning the parents of a vertex is cheaper than that for getting its children. In our implementation, the former involves only dereferencing a pointer and checking that the referenced value is not 0. Furthermore, this check only needs to be performed for coalescence vertices. The latter check requires a single dereference operation as well, but is followed by an arithmetic comparison with a variable and must be performed for every vertex. Performance gains are negligible for small samples, but accumulate as ancestries grow. The complete MMN procedure is given in algorithm 5. Storing traversed edges and associated chunks before applying algorithm 4 (line 7) is wasteful but necessary to exploit SIMD parallelism. Line 14 assumes that markers are encoded right-to-left within chunks (e.g., the sequence 111000 is stored as 000111). On modern architectures, `trailing_zeros` essentially corresponds to the `TZCNT` or `BSF` instruction and is consequently very fast. `idxtopos` (line 17) is the function that computes the position of a marker from its index; see subsection 3.4.2 for details. The chunk size of 1 byte (8 bits), assumed throughout algorithm 5, is a compile-time constant for performance reasons. Although we believe it to be the best choice in most circumstances, it can easily be configured via Julia's preferences mechanism.

### 3.2.2 Breakpoint

Once an inconsistent marker has been found, a sequence of RR events is sampled in a way that reduces the local number of mutations down to one. The RR procedure itself is discussed in details in subsection 3.2.3. For now, we are concerned with the distribution of breakpoints.

Standard theory models the occurrence of recombination events as a Poisson process, which leads to an elementary algorithm for sequential simulators revolving around sampling exponential interarrival times. Things are different when inferring ARGs. Assuming a strong mutation model such as the ISM greatly restricts the number and position of breakpoints. Let $b$ and $m$ be the positions of the last breakpoint and the next inconsistent marker, respectively. Since, under the ISM, some events must be positioned in this interval, it would not make sense for the distance between $b$ and the next breakpoint to follow an exponential distribution. A standard result about PPPs states that events are uniformly distributed in an interval conditional on the number of events in that interval. From this point of view, the correct way of sampling breakpoints in $[b, m)$ would be to first draw a Poisson-distributed number of events and then

---

**Algorithm 5:** MMN

---

1  Set markeridx $\leftarrow 1$

2  **while** markeridx $< \#\,markers$ **do**

3    Empty edges, chunks_src, chunks_dst and mutationedges

4    Set chunkidx to the index of the chunk containing marker markeridx

5    Compute $\omega$, the ancestral interval of unmasked markers in the chunk

6    Compute $m_\omega$ to mask irrelevant markers to 0

7    Fill edges, chunks_src and chunks_dst with edges and adjcent chunks by traversing the marginal graph
     associated with $\omega$ bottom-up, stopping early when a chunk masked with $m_\omega$ is equal to zero.

8    **for** k $\in 1, \ldots, |$edges$|$ **do**
       // Lines 2 and 3 of algorithm 4 (in-place)

9      Set chunks_src[k] $\leftarrow$ chunks_src[k] $\oplus$ chunks_dst[k]

10     Set chunks_src[k] $\leftarrow$ chunks_src[k] $\otimes m_\omega$

11   **for** i $\in 1, \ldots, |$edges$|$ **do**

12     Set acc $\leftarrow 0$

13     **while** *true* **do**

14       Set j $\leftarrow$ `trailing_zeros`(chunks_src[i]) // Hardware operation

15       **if** $j > 8$ **then**

16         break

17       Set pos $\leftarrow$ `idxtopos`(*8(*chunkidx- *1)* + acc + j)

18       **if** pos $\in$ `ancestral_interval`(e) **then**

19         `push!`(mutationedges[acc + j], e)

20       Update chunks_src[i] $\leftarrow$ chunks_src[i] $\gg j$

21       Update acc $\leftarrow$ acc $\gg j$

22   Set mutationidx_chunk to the index of the first entry of mutationedges with cardinality strictly greater than
     one. If there is no such entry, mutationidx_chunk $\leftarrow 0$.

23   **if** mutationidx_chunk $> 0$ **then**
       // ARG is not consistent; return index of the next inconsistent marker &
           associated mutation edges

24     Return $8($chunkidx $- 1) +$ mutationidx_chunk, mutationedges[mutationidx_chunk]

25   Set markeridx $\leftarrow 8$chunkidx $+ 1$

   // ARG is consistent

26 Return $0, \varnothing$

---

position them uniformly in $[b, m)$. Unfortunately, sampling a sequence of events that would yield a tree consistent with the marker at $m$ having a predetermined length is prohibitively difficult unless we are ready to either engage in additional time-consuming computation or exclusively locate recombination events on *derived edges*, that is those edges whose downstream vertex (the one closer to the leaves) has the derived allele at the focal marker. As described in subsection 3.2.3, the maximum number of mutations eliminated by a recombination event occurring on a *wild edge*, that is, an edge whose downstream incident vertex is associated with the wild allele for the marker under consideration, is the number of mutations in the marker's marginal tree minus one. In other words, depending on the ARG's topology, a constrained RR event might remove anywhere from a single mutation to all of them.

An additional complication is that while sampling breakpoints in $[b, m)$ leads to a reduction in the number of mutations for the marker at $m$ assuming correct locations for RR events, there is no guarantee that doing so will not create mutation edges for markers to the left of $m$. Since we are inferring ARGs from left to right, increasing the number of mutation edges for markers to the right of $m$ is not a problem. It would generally be desirable for that number to go down, but an algorithm effectively achieving this goal would likely not scale. In any case, additional mutation edges created for markers to the right will be dealt with in subsequent iterations. Creating mutations for markers that have already been processed is more problematic. While we could backtrack to deal with those, this would negatively affect performance and necessitate sampling additional events. Instead, Moonshine imposes additional constraints on

breakpoints' distribution. Let $m'$ be the position of the first marker to the left of $m$. An event positioned in $[m', m)$ cannot create mutation edges for markers to the left since any associated edge is not ancestral for those markers.

Restricting the support of breakpoints to the interval between the next inconsistent marker and the one directly to its left is efficient but overly prohibitive. In addition to reducing the support of breakpoints for no reason other than ease of implementation, it diminishes our algorithm's ability to sample type 2 recombination events. Algorithms based on a first-order Markov approximation of the recombination process, such as SMC and SMC', take a rather drastic approach by simply ignoring any event that might have happened anywhere except on the current marginal tree. Other algorithms allow recoalescence events to happen further along the sequences while still restricting the position of recombination events. The main selling points of these algorithms are their computational efficiency and relative simplicity. However, since recombination events are precluded from happening in non-ancestral material, they are inherently incapable of sampling type 2 recombinations. In order to simulate these, a sequential algorithm must be able to go back in space, so to speak, while simulating the recombination process: some events must happen to the left of the preceding one. This class of algorithm includes samplers such as SC and, by extension, SC-sample. Their strategy is to sample a recombination event on either the coalescence branch of the previous RR event or one located upstream. This process is repeated until recoalescence with the current marginal tree occurs. ARGinfer is another algorithm able to sample type 2 recombination. Being an MCMC sampler, it does so via two proposals: "adding a new recombination to a lineage" and "resampling the breakpoint of a recombination event" ([23]).

These algorithms have different approaches to sampling type 2 events. One approach, used by SC, aims to stay closer to the coalescent process by sampling recombinations according to a complex distribution. The other, used by ARGinfer, uses a less sophisticated distribution and relies on the properties of its underlying MCMC sampler for statistical correctness. It is apparently able to propose a move faster than SC can sample a recombination event. However, many of those moves may be rejected before one is finally accepted, while every event sampled by SC is actually integrated into the ancestry Both approaches have their merits, and it would be risky to make a unilateral statement as to which is the best.

Our sampler shares many characteristics of SC and ARGinfer, but its approach to type 2 events is different. Since we are processing haplotypes from left to right and only sample events resulting in a reduction of the marginal number of mutation edges, breakpoints associated with type 2 events have to be positioned to the left of not only $m$ but of $m'$ as well. This begs the question: what are the conditions under which a recombination event positioned at $b \leq b' < m'$ does not increase the number of mutation edges in $[b, m)$? Let $e_r$, $e_c$ be the recombination and recoalescence edges respectively, $v_r$, $v_c$ the associated child vertices, and $h_{v_r}^k$, $h_{v_c}^k$ the status of the $k^{\text{th}}$ marker of the associated haplotypes. An RR event positioned at breakpoint $b'$ will not result in the creation of additional mutation edges in $[b, m)$ as long as for every marker $k$ in $[b', m)$, one of the following conditions is met:

1. $h_{v_r}^k = h_{v_c}^k$;
2. $e_r$ is not ancestral for marker $k$;
3. $e_c$ is not ancestral for marker $k$ and the status of the first ancestral vertex upstream $v_c$ is $h_{v_r}^k$.

Among these conditions, item 2 allows for type 2 events. Since we are working from left to right, support for those breakpoints comes with the potential of sampling type 4 events as well. Consequently, they must be explicitly discarded.

Assuming the lower limit of the support for the $i^{\text{th}}$ breakpoint $b_i'$ is known, we need to decide on a distribution for the breakpoint itself. Since the support is bounded, the two most natural choices are the truncated exponential and uniform distributions. The main advantage of the former is its resemblance to the untruncated exponential distribution, which is the mathematically correct choice assuming a PPP and an unknown number of recombination events in $[b, m)$. For minimum departure from this model, the $i^{\text{th}}$ breakpoint $b_i$ would be supported on $[\max\{b_i', b_{i-1}\}, m)$ where $b_0 = 0$. In addition to precluding type 2 events, this distribution has the drawback of being difficult to deal with from a numerical standpoint. Each draw reduces the support of the next, making standard approaches to random generation such as rejection sampling and inverse transform less efficient and/or more prone to instability due to divisions and the use of transcendental functions. Consequently, Moonshine uses the simpler uniform distribution. Each breakpoint is supported on the full interval $[b_i', m)$ which avoids numerical instabilities and allows for type 2 events. This is akin to trying to stay as close as possible to the PPP distribution assuming (erroneously) a known number of events.

### 3.2.3  Recombination and Recoalescence

When processing sequences from left to right, the procedure described in the previous section allows to efficiently find the next site incompatible with the current ARG, that is, the site with two or more mutation edges. It is straightforward
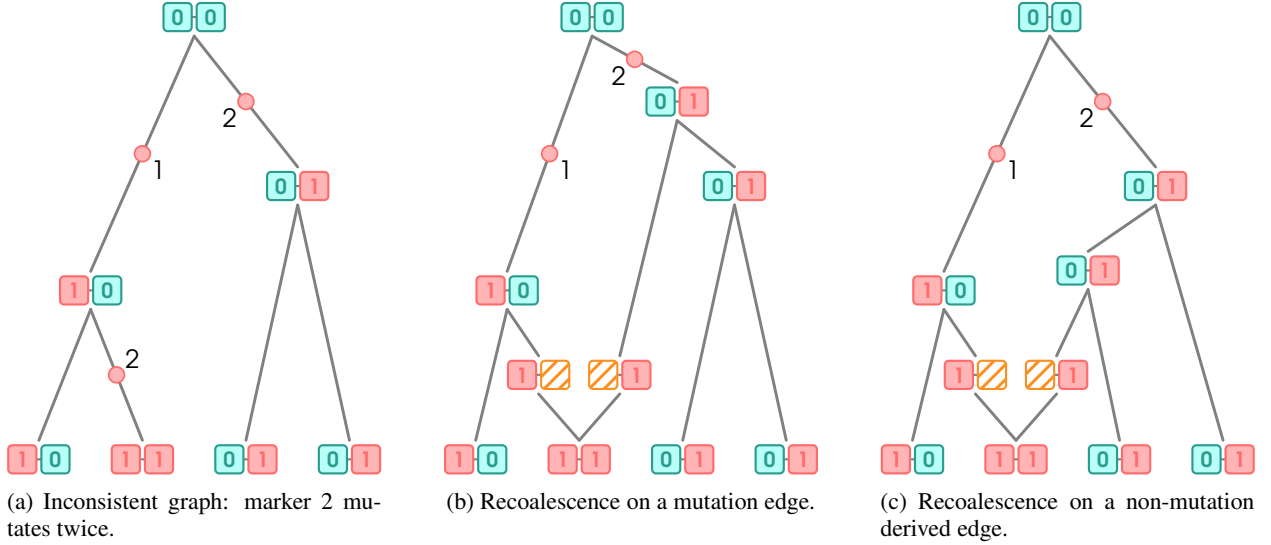
(a) Inconsistent graph: marker 2 mutates twice.

(b) Recoalescence on a mutation edge.

(c) Recoalescence on a non-mutation derived edge.

Figure 4: Derived recombination followed by a recoalescence event on a derived edge. Original graph is represented in fig. 4a. Figures 4b to 4c are example of the two possible types of recoalescence events.

to keep track of these by storing them in a list, for instance. The next step is then to add recombination events to the ARG in a way that renders it consistent with the focal site.

Let $M$ be the set of mutation edges. The coalescence of two elements of $M$ results in a net diminution of the number of mutation edges by 1. The reason is as follows: since the two coalescing edges have the same allele at the focal site (namely the derived one), neither will be a mutation edge anymore after coalescence. However, one edge upstream of the coalescence vertex will ultimately coalesce with a wild edge and therefore become a new mutation edge, resulting in the aforementioned net diminution. To conclude the argument, note that it is not possible for a coalescence event to decrease the number of mutation edges by more than 1. For that to be possible, one of the coalescing edges would need to have a sibling that is itself a mutation edge. However, in that case, the coalescing edge would not be a mutation edge. We call recombination events located on a derived edge and the recoalescence event that follows *derived RR event*. A configuration resulting from such a pair of events is given in fig. 4b. The previous argument actually entails an even stronger conclusion: there is no requirement for the edge on which recoalescence occurs to be a mutation branch. As illustrated in fig. 4c, merely being derived is sufficient.

The type of derived events we just described is said to be *constrained* because their locations and positions are supported on a subset of the ARG branches and sequences, respectively. Ideally, this subset should encompass as many branches as possible, making the distribution closer to the true (i.e., PPP) one. It is in fact possible to extend the set of candidate branches without compromising the reduction of the number of mutations. For one thing, recombination events do not have to be located on a mutation edge. Let $u - v$ be such an edge with $v$ being the downstream vertex. Since $v$ is derived, so is every vertex located downstream in the current marginal tree. Let $E_v$ be the set of edges downstream of $v$ in that tree. The number of mutations can be reduced by one by taking a subset of $E_v$ separating $v$ from the ARG's leaves and generating an RR event on each edge, provided that each recoalescence events is located on derived edges outside $E_v$. If we want to avoid inflating the number of recombination events, we can impose the additional restriction $|E_v| = 1$ and sample a recombination event on $u - v$ if no such edge separator satisfies this constraint. This is easily achieved in practice by looking for sequences

$$u - v, v - w_1, w_1 - w_2, \ldots$$

where each edge is marginally without sibling. In particular, this is the case if $v, w_1, w_2, \ldots$ are recombination vertices, although it is not a necessary condition, as illustrated by fig. 5b. Considering these edges as possible locations for recombination events allows for a more even distribution, which, in addition to making the approximation to the true distribution, contributes to the elimination of numerical errors associated with a large number of such events relative to the total branch length. Sequential algorithms come with something of a numerical curse: as the ARG grows, so does the number of short branches resulting from events being sampled close to a branch's incident vertices. Soon enough, numerical instabilities arise when a short mutation edge is encountered. Fortunately, as we have seen earlier, the remedy is simple to implement.

(a) Inconsistent graph: marker 2 mutates twice.

(b) Recombination event on one of the child edge of a mutation edge, its sibling being non-ancestral for marker 2.

(c) Recombiation event on a mutation edge followed by recoalescence on a non-ancestral branch downstream of another mutation edge.
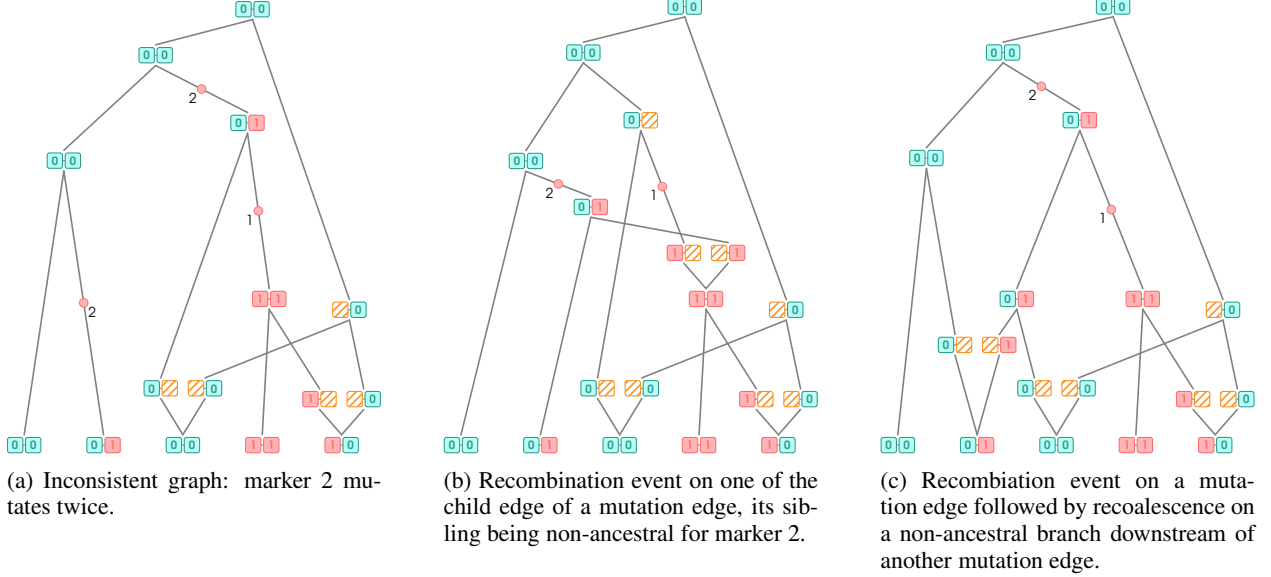
Figure 5: Derived recombination event located on an edge downstream of a mutation edge. Reduction in the number of mutation is achieved since the recombination edge has no sibling marginally.

We go even further by applying a similar analysis to the location of recoalescence events. As noted by [10], these do not need to be supported on derived edges. A recoalescence with a non-ancestral edge downstream of a derived edge does not increase the number of mutations, a scenario illustrated in fig. 5c. Since it is impossible for a wild edge to lead to a derived one, the set of admissible locations for recombination events can be described in rather simple terms: coalescence of a derived edge with a mutation edge or any of its downstream edges does not increase the number of mutations.

It should be clear by now that we can make any position consistent by sampling at most $|M| - 1$ derived recombination events. This is good, but we can actually do better, at least some of the time. We refer to recombination events located on a wild edge and the subsequent recoalescence event as *wild RR event*. Despite not occurring on a mutation edge, a wild recombination event can effectively reduce the number of mutations when the following two conditions are met:

1. The recombination edge's brother is a mutation edge;
2. Their uncle is a mutation edge.

In such a configuration, the recombination event turns both derived vertices, namely the sibling and uncle of the recombination vertex, into siblings with respect to the current marginal tree, reducing the number of mutations by one. As with the derived case, the wild recombination event does not need to occur on the sibling of the mutation edge; the number of mutations will be reduced as long as it is located downstream somewhere along a "chain" of edges without marginal siblings. Just like in the derived case, coalescence can occur with non-ancestral edges as long as they lead to a wild branch; this is illustrated in fig. 6. Failing to meet this condition would result in the creation of a new mutation event.

When it comes to the number of mutations eliminated, wild RR events clearly have the upper hand. As we just discussed, a wild recombination event flips the allelic state of the mutation edge's parental vertex to derived. Consequently, a reduction in the number of mutations can propagate upward as long as the uncles encountered along the way in the current marginal tree are also mutation edges. In fact, given an appropriate topology, a single mutation event can reduce an arbitrary number of mutations to a single one. A scenario where a single coalescence/recombination event leads to the elimination of two mutations is illustrated in fig. 7. On the contrary, as discussed before, the number of mutations eliminated by a derived RR event is limited to one. This is because the mutation is eliminated by the recoalescence rather than the recombination event. The state of the marker of the parental vertex of the mutation edge on which the recombination event occurred remains unchanged, inhibiting the propagation phenomenon observed in the wild case.

That being said, although `Moonshine` has the capability of inferring ARGs by sampling mutation reducing events exclusively, it does not actively seek to maximize mutation reduction locally. It does not, for instance, give priority to wild events over derived ones. As our main objective is to remain close to the CWR, only branch length is taken into
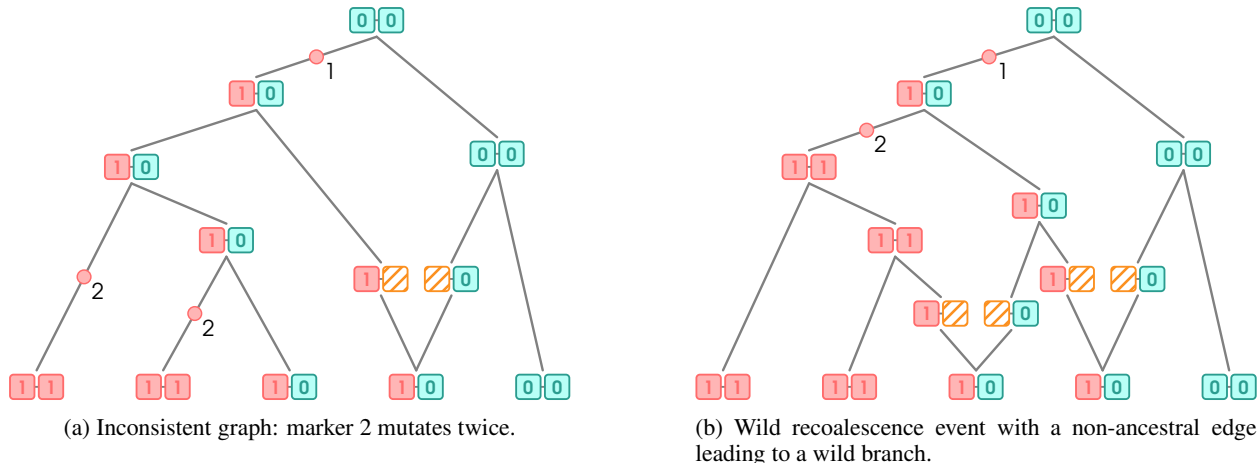
(a) Inconsistent graph: marker 2 mutates twice.

(b) Wild recoalescence event with a non-ancestral edge leading to a wild branch.

Figure 6: Wild RR event leading to a reduction in the total number of mutations with recoalescence on non-ancestral edge.



Figure 7: A single recombination event leads to the elimination of multiple mutations.

consideration when sampling the location of a recombination event. Each edge on which a recombination event would result in a reduction in the number of mutations for the current marker has a probability of being chosen proportional to the difference in latitude of its incident vertices, a sampling strategy described earlier as "constrained". For reasons discussed earlier, the event's location on the branch is not sampled uniformly but rather assumed to be distributed as the same location-scale family associated with the Beta distribution as the one used for the unconstrained case.

Determining the location of a derived recoalescence event is fairly straightforward. According to standard theory ([7]), coalescence occurs at a unit rate with each admissible edge. We simply need to list all available locations and pick one uniformly at random. The set of possible recoalescence edges $E_R$ is readily established and does not require graph traversal. The minimum latitude is either that of the recombination event or the smallest latitude among the destination vertices of $E_R$, whichever is greater. Each edge in $E_R$ is associated with a probability proportional to its length minus any section outside of admissible latitudes. Once the recoalescence branch is determined, a location is sampled as usual on its admissible portion.

Simulating the location of a wild recoalescence event is more involved due to the semi-infinite nature of its support. Its latitude is distributed as the time of the first event of a non-homogeneous Poisson process. A common approach in one-dimensional scenarios such as ours is the *time-scale transformation* method, which is a form of inverse transform sampling and requires computation of the inverse of the integrated intensity function, also known as the *cumulative intensity function*. The main issue stems from our decision not to track the number of live edges by latitude. This improves general performance and reduces memory usage, but makes evaluation of the intensity function extremely time-consuming, as it requires graph traversal. We tackle this issue using numerical integration. Although somewhat variable, the intensity function is piecewise constant and therefore a very good candidate for quadrature. We use a logarithmic grid of latitudes to account for the generally decreasing complexity of the ARG topology as height increases. The intensity function is evaluated at quadrature nodes in a single partial traversal of the ARG, reducing overhead to the minimum. By default, our algorithm uses a grid of compile-time constant size 25, but this number can be tuned to balance precision with performance using the preference mechanism.
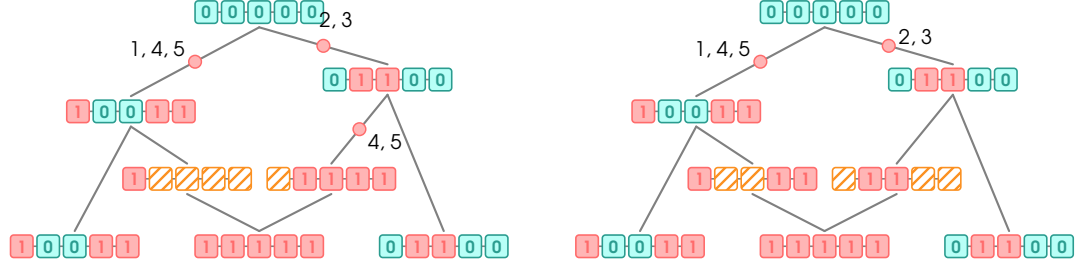
16

Figure 8: Multiple crossing over.

### 3.2.4 Multiple Crossing Over

In addition to RR events, `Moonshine` can sample multiple crossing over (MCO) events constrained to reduce the number of mutations. These events arise in the following scenario: assume that the branch on which a recombination event has been sampled is the right parental edge of a recombination vertex. Assume further that a recoalescence with the sibling of the other parental edge's branch would reduce the number of mutations. When both of these conditions are met, RR events are unnecessary; the same effect can be achieved by modifying the ancestral intervals of the parental edges of the recombination vertex. For this to work, it is necessary to track the partition induced by recombination events. In an `ARG`, every recombination vertex is associated with two sets of intervals, one for each parental edge. We call these sets *recombination masks* and denote the mask associated with recombination $k$ by $m_k = \{m_k^l, m_k^r\}$. Initially, before any MCO event involving a specific recombination vertex, the mask is rather simple. Let $b$ be the position of the associated recombination event.

$$m_k^l = [0, b)$$
$$m_k^r = [b, \infty).$$

A MCO event at position $b' > b$ transforms those sets as follows:

$$m_k^l = [0, b) \cup [b', \infty)$$
$$m_k^r = [b, \infty) \cap [0, b') = [b, b').$$

A subsequent event at position $b'' > b'$ would yield

$$m_k^l = [0, b) \cup ([b', \infty) \cap [0, b'')) = [0, b) \cup [b', b'')$$
$$m_k^r = [b, b') \cup [b'', \infty).$$

In general, the correct mask can be computed by intersecting the rightmost interval (with respect to the right endpoint) with $[0, b')$ and taking the union with the other interval and $[b', \infty)$. Although explicit storage is not necessary, coalescence vertices can be thought of as being associated with the mask $[0, \infty)$. Recombination masks are used to compute the ancestral intervals of parental edges in the following way: an edge's ancestral interval is equal to the intersection of its recombination mask with the union of its children's intervals.

Conditions under which a MCO event can occur are very restrictive. Consequently, we expect detection to be underpowered. This might be improved in the future, for example, by allowing the user to bias the sampler towards these events. A MCO is illustrated in fig. 8.

### 3.3 ARG Update

RR and MCO events both modify sequences and ancestral intervals associated with vertices and edges upstream. The affected elements must be updated immediately to ensure the soundness of the remainder of the procedure, which can be achieved by traversing the ARG from the recombination and recoalescence edges toward the root. This seemingly simple operation is, however, computationally demanding. Our first, rather naive implementation was a major bottleneck of the constrained recombination algorithm. To make it as efficient as possible, we limit the update procedure to the elements affected by the event. Recombination events are mostly local, meaning they can only affect vertices and edges located upstream. Consequently, any element not upstream of either the recombination or recoalescence edge can be ignored when updating. In fact, an element must be located upstream of a modified element to be modified itself. This means that the number of updates can be further reduced by keeping track of the state of every vertex and edge before they are updated and stopping when a match between original and updated versions is detected. The algorithm terminates early if both the sequence and the ancestral interval associated with an edge are left unchanged.

17

The early termination strategy described above dramatically decreases the time dedicated to ARG update. Indeed, coalescences with vertices left untouched limit the spread of changes, often well below the root. It might be conceived that the additional costs associated with storing information about elements before updating them would outweigh the benefits of reducing the number of updated elements. It turns out that a very significant number of recombination events are very local in nature, to such a degree that we have yet to find the point of diminishing return for reasonably large samples. That is not to say, however, that the procedure cannot be improved further. We were able to squeeze even more performance out of it through hashing. Instead of making a copy of the current sequence and ancestral interval before update, we simply compute a hash value for each of those, which we promptly hash together. The procedure is terminated early if the hash of the updated sequence-ancestral interval pair is equal to the original. Since hash functions are not injective, there is a possibility that different pairs may have the same hash, an event known as *collision*. Fortunately, this has not been a problem in practice. To put any doubt to rest, the method `validate` can be applied to the final product of the ARG inference process to ensure that it is exempt, among other things, of the inconsistencies that would emerge from collisions. The complete procedure is presented in algorithm 6. Note that since recombination vertices have a unique child, their associated sequence can be a reference to their child's. Consequently, line 5 can be performed without additional allocation, dramatically cutting down on memory usage. The early termination strategy is implemented by line 16.

---

**Algorithm 6:** ARG Update

---

1  Set stack ← [starting_vertex]
2  **while** stack *is not empty* **do**
3  | Set v ← pop!(stack)
4  | **if** v *is a recombination vertex* **then**
5  | | Set the sequence of v to that of its child
6  | | Set the ancestral interval of its parent edges to the intersection of its child's with the appropriate ancestral mask
7  | | push!(stack, parents(v))
8  | **else**
9  | | Set h ← hash(haplotype(v))
10 | | Set the sequence of v to the conjunction of its children's haplotypes
11 | | **if** v *is the root* **then**
12 | | | Continue
13 | | Set e to the parent edge of v
14 | | Update h ← hash(h, ancestral_interval(e))
15 | | Set the ancestral interval of e to the union of the ancestral intervals of v's child edges
16 | | **if** h ≠ hash(hash(haplotype(v)), ancestral_interval(e)) **then**
17 | | | push!(stack, parent(v))

---

### 3.4  Technical Considerations

#### 3.4.1  Markov Approximation

The procedure described in this section is exact in the following sense: at any step, recombination and recoalescence edges are free to be sampled anywhere on $G$ subject only to the mutation number reduction constraint. It is more similar to the original Wiuf-Hein sequential algorithm than any of its Markovian approximations, such as the SMC or SMC'. This makes our sampler more realistic, but it comes with a performance penalty: as RR events are incorporated into the graph, the number of edges that need to be considered for the next events increases. Consequently, we should expect the number of operations performed by our algorithm to grow as it progresses along sequences. One way to alleviate this computational burden is to make the recombination process artificially Markovian. The procedure we just presented allows for such an approximation, although it was omitted from the description for simplicity. Both restricted and unrestricted RR events may be executed inside a window moving across the sequences, the width of which can be selected by the user. A width of 0 corresponds to a first-order Markovian approximation akin to the SMC, while an infinite width means no approximation at all. In general, assuming distances in base pairs (bp), specifying a width of $w$ for an RR event occuring at position $p$ has the effect of excluding any edge $e$ such that

$$\text{ancestral\_interval}(e) \cap [p - w, p + w] = \varnothing.$$
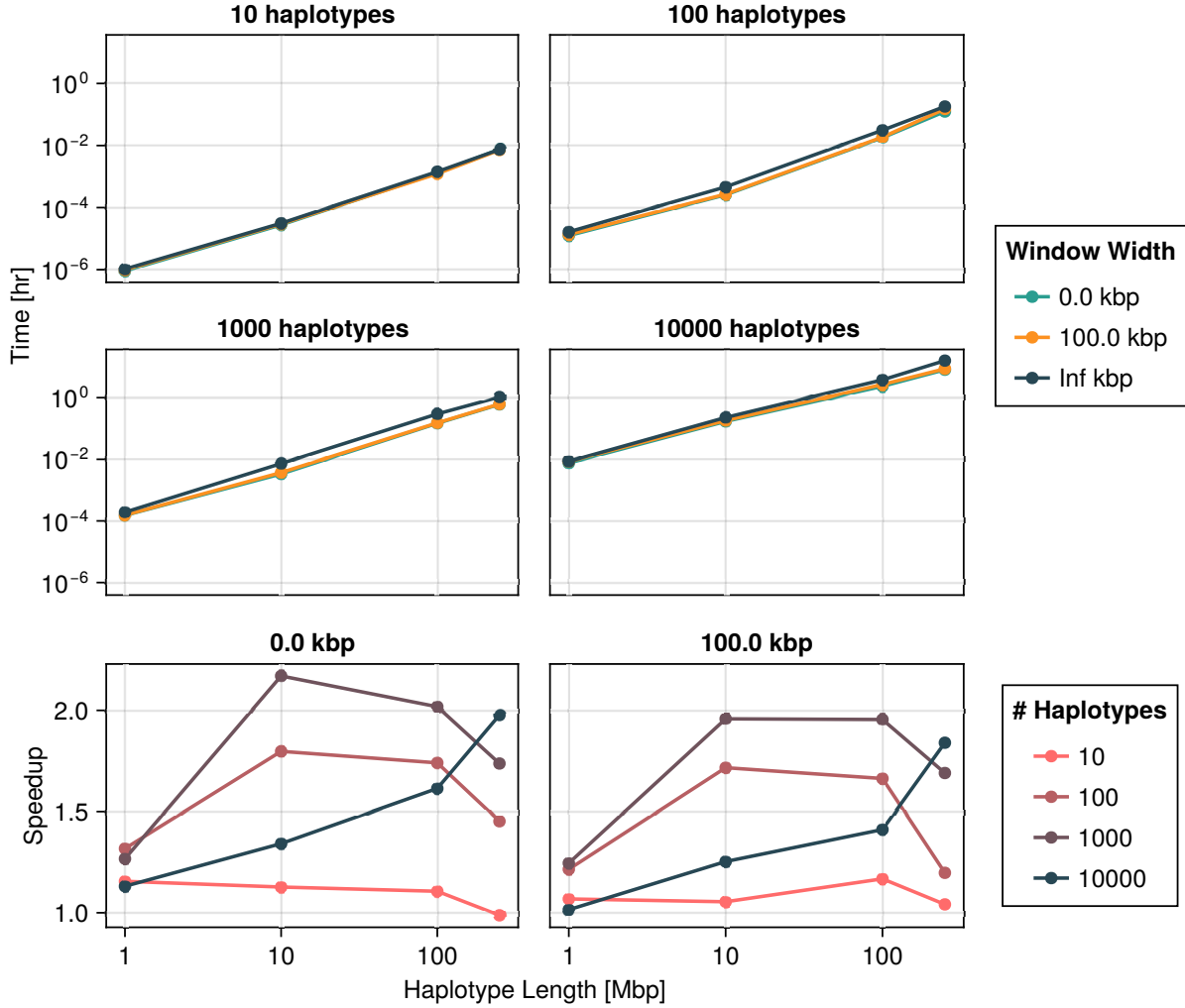
18

Figure 9: Time required to sample a single consistent ARG as a function of the number of haplotypes and sequence length, including time spent inferring the initial tree. The speedup provided by various finite window widths is depicted in the lower figure. Technical details are discussed in subsection 5.1.

from the set of candidate recombination and recoalescence edges. We emphasize that the window is *centered* on $p$, as our sampler is designed for the general task of rendering ARGs consistent rather than building them from a tree in a left-to-right sweep. Moreover, when sampling the recoalescence latitude of an unconstrained event, ignored edges are not taken into account for the computation of the rate of the recombination latitude. Figure 9 shows that the performance impact of reducing the window's width is significant. For some combinations of sample size and haplotype length, a window width of 0 can lead to over twice as fast computation. Interestingly, the speedup is similar for a width of 100 kbp, suggesting that even relatively modest approximations can lead to substantial reduction in computational resources. Irrespective of window size, fig. 9 suggests that computation time is exponential in both the number of sampled haplotypes and markers.

Similarly, fig. 10 suggests exponential growth in memory usage. It appears, however, that reducing the window size increases the size of the resulting ARG. This is largely explained by the increasing number of recombination events sampled as a result. As shown in fig. 11, this number tends to increase as the window width diminishes. Part of this phenomenon might be explained by increased flexibility. Larger window sizes increase the number of candidate edges for recombination and recoalescence events, allowing for constrained coalescence of more similar haplotypes and, ultimately, more parsimonious ancestries. This behavior is all the more interesting as it is entirely spontaneous;
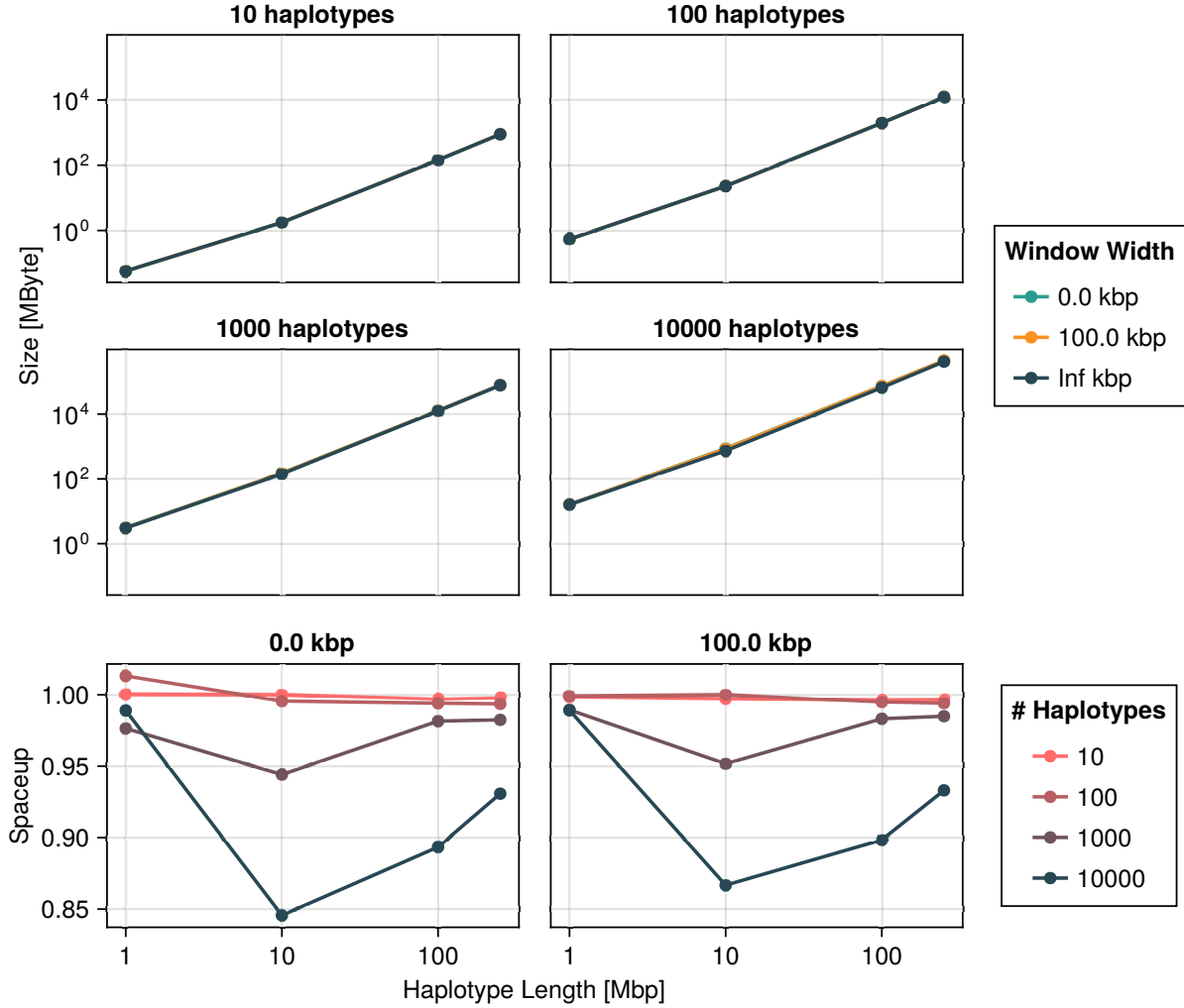
19

Figure 10: Size of the object containing the inferred ancestral recombination graph as a function of the number of haplotypes and sequence length. The spaceup is the ratio of the size for infinite window width versus that for the indicated width. Technical details are discussed in subsection 5.1. We emphasize that the total amount of memory required by the sampling procedure exceeds quantities reported here.

markers to the right of the recombination breakpoint are not taken into account in the recoalescence edge sampling procedure.

### 3.4.2 Markers and positions

In addition to the ARG update procedure described at the beginning of this section, another major bottleneck in our method was, somewhat unexpectedly, a function called `postoidx`. It is the inverse of another function called `idxtopos`, which is designed solely to return the position of a given marker, indexed from the leftmost one. The term "inverse" in that context is used loosely, as only a few positions are associated with a sampled marker. For that reason, we define `postoidx` formally as a pseudoinverse of `idxtopos`:

$$\texttt{postoidx}(p) = \sup\{i \in 1, \ldots, s : \texttt{idxtopos}(i) \leq p\}.$$

Since it is often necessary to mask sequences with respect to an ancestral interval, this function is used extensively within `Moonshine`'s recoalescence and ARG update procedures. The endpoints of any given interval rarely correspond to the position of a marker. This is compounded by the fact that what we refer to as "ancestral interval" is, in reality, the union of multiple disjoint intervals. Overall, a simple linear search is insufficient; a more efficient approach is
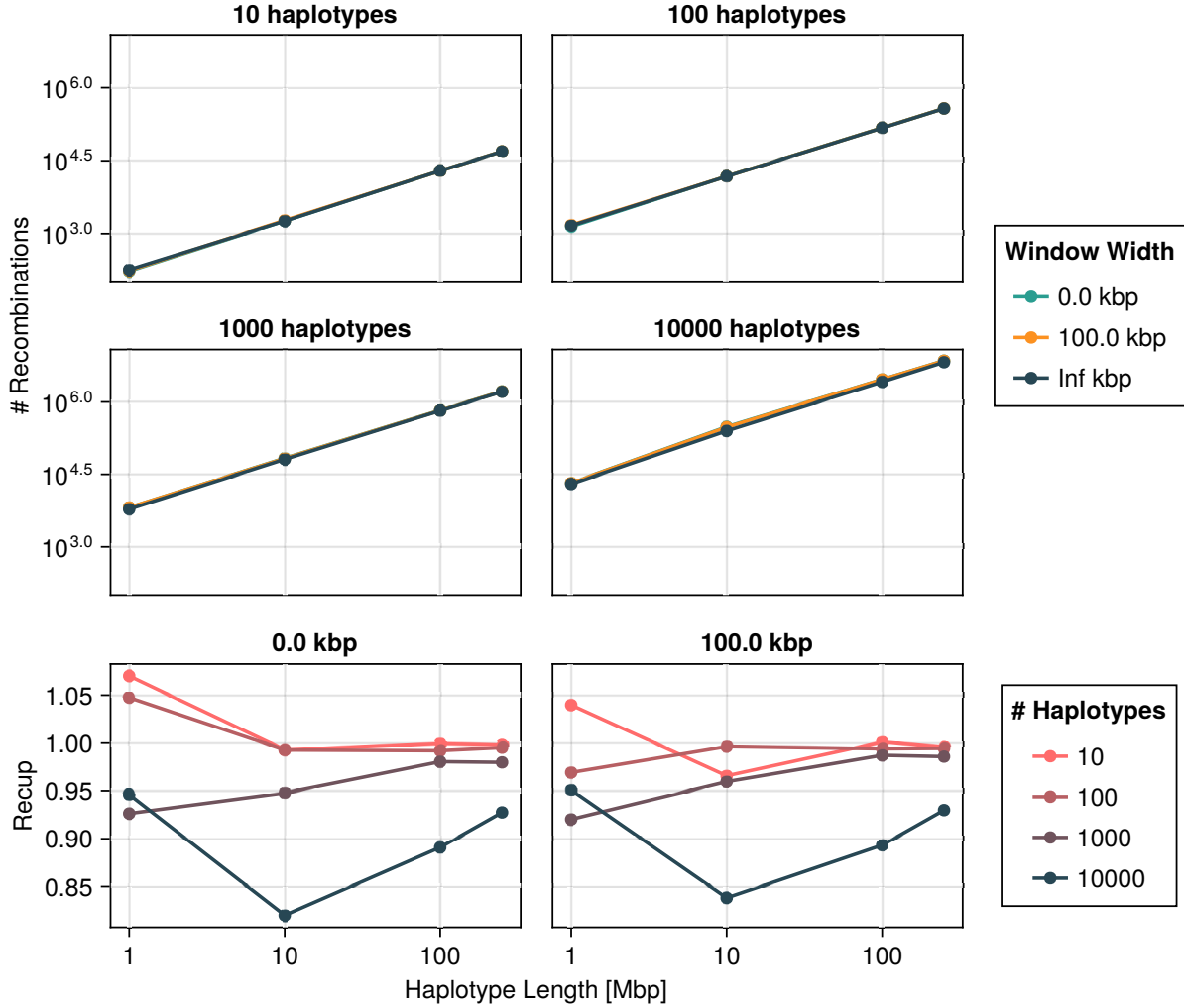
Figure 11: Number of recombination events sampled by the ARG inference procedure as a function of the number of haplotypes and sequence length. The recup is the ratio of that number for an infinite window versus that for the indicated width. Technical details are discussed in subsection 5.1.

called for. Binary search yields considerable gains, but we can do better. Memoization is a standard strategy for solving similar problems, but it is detrimental to our use case. The cost of computing `idxtopos` via binary search is small, even in comparison with a lookup in a hash table. Besides, positions are encoded as double-precision floating-point values, which limit their usefulness as keys of an associative structure. Perhaps it would be possible to use locality-sensitive hashing to circumvent this issue. We did not explore this avenue, though, because the approach currently implemented, although relatively simple, proved satisfactory. It is essentially a slightly optimized version of interpolation-sequential search ([35]). We take advantage of the fact that the markers' position vector is not modified by ARG inference. Instead of directly computing the inverse of `idxtopos`, we do so on a first-order approximation obtained by least-squares fitting, which is straightforward. Even though computing such an approximation is time-consuming relative to an iteration of bisection search, it only has to be done once, making the overhead negligible. Pseudocode for `postoidx` is given in algorithm 7.

While algorithm 7 is relatively simple, it takes advantage of three features of the positions vector: its static nature, its monotonicity, and the approximate regularity of the distance between markers. A more efficient procedure could likely be designed, but algorithm 7 is at least efficient enough to eliminate the `postoidx` bottleneck.

---

**Algorithm 7:** postoidx

---

// Assume idxtopos(i) $\approx$ ai $+$ b

**1 Function** postoidx(p)

**2**     Set i $\leftarrow \frac{p-b}{a}$, l $\leftarrow 1$, r $\leftarrow$ s

**3**     Clamp i in $[1, s]$

**4**     Update l or r through a single binary search iteration

**5**     Perform interpolation-sequential search with i, l and r

---

### 3.4.3 Memory Allocation

Considerable effort has been devoted to using memory as efficiently as possible in the computationally demanding sections of Moonshine. In addition to the obvious benefits in terms of reduced memory requirements, we found dynamic allocation to be a major performance bottleneck in itself. Although Julia is a high-level programming language, fine-grained memory management is straightforward. It is entirely possible to manage allocations directly. Part of the C standard library, comprising functions such as malloc, calloc and free, is exposed to the user via the Libc module included in Julia's standard library. Our approach, however, is slightly higher level: we make extensive use of slab allocation ([36]) through Bumper.jl ([37]). This package allows us to handle raw pointers when needed while efficiently managing the underlying memory pool automatically. This allows for nearly effortless tight memory management and dramatically reduces the number of allocations needed in core methods. Performance-critical methods accept the buffer keyword argument through which the caller can pass a Bumper.jl buffer. This is how the memory is allocated by the build! method for inferring ARG for instance. The same buffer can be passed around to different function calls for maximum memory usage efficiency. This also makes it extremely easy to seamlessly integrate Moonshine methods in a workflow that already takes advantage of Bumper.jl facilities.

### 3.5 Complete Algorithm

The complete ARG inference procedure is given in algorithm 8.

## 4 Conclusion

We have described a fast and efficient algorithm for ancestral recombination graph inference and the broad outlines of its Julia implementation via the Moonshine package. In addition to exact sampling, an approximate scheme based on user-defined genomic windows is available, which reduces inference time. Inference is grounded in a restricted version of the coalescent-with-recombination distribution rather than heuristics. We have also presented a flexible algorithm for coalescent tree construction that can accommodate a wide range of metrics on haplotypes, exactly or approximately, while avoiding numerical instabilities.

As of version 0.3.3, these algorithms are for the most part mature. No significant leap in inference time or memory usage is to be expected in the near future. Our goal for version 1.0.0 is to improve real-world usability. Specifically, the API reference needs to be completed to make implementing alternative models a more pleasant experience. General documentation and guides also need to be improved for the general user. Both are publicly available at https://moonshine.patrickfournier.ca.

As of today, version 0.3.4 and later are compatible with files encoded in the variant call format (VCF) through the VCFTools.jl package ([38]). We plan to support FASTA/FASTQ and sequence/binary alignment map format (SAM/BAM) via the FASTX.jl and XAM.jl packages, respectively. Along with compatibility with these input formats, we plan to support missing markers by treating them as non-ancestral material, which is the approach taken by [39, 28, 40].

Moonshine is available on Julia's general registry and can be easily and quickly obtained using standard facilities. In addition, we plan to package it in two application containers. One will include the Pluto notebook ([41]) package and graphical utilities such as Makie.jl ([42]). Its goal is to reduce the burden associated with ARG inference for practitioners. The other container will be more minimal, including only the minimum required to deploy Moonshine via an orchestration system such as kubernetes. Our hope is to facilitate the creation of high-performance computation clusters in small-to-medium environments.

---

**Algorithm 8:** ARG Inference

---

1 Set markeridx $\leftarrow 1$, live_edges $\leftarrow \varnothing$
2 Update markeridx and live_edges via algorithm 5

3 **while** $markeridx > 0$ **do**
4    **while** $|\text{live\_edges}| > 1$ **do**
5       Sample $e_1$ & $e_2$ uniformly from live_edges and remove them

      `// Consider the possibility of a wild recombination`
6       Set markerpos $\leftarrow$ `idxtopos(`markeridx`)`
7       **if** `parent(src(`$e_1$`),` markerpos`) = src(`$e_2$`)` **then**
8          Set redge $\leftarrow$ `src(`$e_1$`) −` `sibling(dst(`$e_1$`),` markerpos`)`
9          Set vupdate $\leftarrow$ `src(`$e_2$`)`
10       **else if** `parent(src(`$e_2$`),` markerpos`) = src(`$e_1$`)` **then**
11          Set redge $\leftarrow$ `src(`$e_2$`) −` `sibling(dst(`$e_2$`),` markerpos`)`
12          Set vupdate $\leftarrow$ `src(`$e_1$`)`

13    Sample a recombination edge redge with probability proportional to its length among all admissible edges (see subsection 3.2.3 for admissibility criteria)
14    Sample a recombination latitude rlat betwen `dst(`redge`)` and `src(`redge`)` from a position-scale family associated with the Beta(2, 2) distribution

15    **if** *Recombination is on a wild edge* **then**
16       Sample a recoalescence latitude clat distributed as a non-homogeneous Poisson process on $[\text{rlat}, \infty)$ via empirical supremum rejection sampling (see subsection 3.2.3 for complete distribution)
17       Sample a recoalescence edge cedge uniformly among those intersecting latitude clat
18    **else**
19       Fill possible_cedges with derived edges located above redge
20       **for** $k \in 1, \ldots, |\text{possible\_cedges}|$ **do**
21          Set r as to the latitude of the lowest valid recoalescence edge associated with e according to subsection 3.2.3
22          Set minlatitudes $\leftarrow \max\{\text{rlat}, \text{r}\}$
23       Set ubound $\leftarrow \max\{(\texttt{latitude} \circ \texttt{src})(\text{e}) : \text{e} \in \text{possible\_cedges}\}$
24       Set lbound $\leftarrow \min$ minlatitudes
25       Sample recoalescence latitude clat distributed as a non-homogeneous Poisson process on $[\text{lbound}, \text{ubound}]$ via rejection sampling
26       Sample associated recoalescence edge cedge (see subsection 3.2.3 for complete procedure)
27       Set vupdate to the number of vertices plus 2

28    Sample breakpoint conditional on redge, cedge and markeridx (see subsection 3.2.2)

29    **if** `src(`cedge`)` $\in$ `parents(dst(`redge`))` *and* `dst(`redge`)` *is a recombination vertex* **then**
30       Perform an MCO event on vertex `dst(`redge`)`
31    **else**
32       Sample a recombination event on redge at rlat followed by a recoalescence event on cedge at clat at position breakpoint

33    Update ARG via algorithm 6 starting at vupdate
34    Update live_edges
35    Update markeridx and live_edges via algorithm 5

---

# 5 Appendix

## 5.1 Simulations

All simulations were performed with `Moonshine` version 0.3.3 and Julia 1.11.3 on AMD EPYC 9654 processors. Execution time was measured using the `@timed` macro. To reduce variability due to the execution environment and differences between ARGs, each reported measurement is the minimum of three runs, averaged over five graphs. Haplotypes are generated by `msprime` using `StandardCoalescent` and `BinaryMutationModel` as the ancestry and mutation models, respectively. Per-locus mutation and recombination rates were both set to $10^{-8}$. An effective population size of $10^4$ was specified. The calls to `msprime.sim_ancestry` and `msprime.sim_mutations` were made through the `Moonshine` interface to `msprime`. Moonshine 0.3.3 does not implement the recoalescence latitude sampling procedures described in subsection 3.2.3. In particular, its algorithm for wild events is less performant. Consequently, we expect the performance of versions 0.3.6 onwards, which implement the time-scale transformation approach, to be slightly better when using the default grid.

All the code used for the simulation studies, as well as the raw data, is publicly available on Codeberg (`https://codeberg.org/ptrk/moonshine.jl-papers`). For convenient reproducibility, code to execute each simulation presented and produce related figures is grouped in a single Pluto notebook. As the simulations are computationally intensive and require considerable time to complete, it may be more convenient to run them on remote machines. For that purpose, the notebook can be executed as a standalone Julia script. Figure generation steps will be skipped when doing so; they can be run locally later on.

### 5.1.1 Results

Raw simulation results are presented below. The following abbreviations are used in the column names:

- **n**: number of haplotypes in the sample
- **L**: length of sampled haplotypes
- **SeU**: speedup
- **SaU**: spaceup
- **RU**: recup
- **W**: window width
- **Recs**: recombination events

Numbers of markers and recombination events are averaged over the sampled ARGs.

Table 1: Tree Building Time

| n | L [Mbp] | Distance | # Markers | Time | Size [MByte |
|---|---------|----------|-----------|------|-------------|
| 10 | 1 | Hamming | 588.4 | < 1s | < 1 |
| 10 | 1 | LeftM | 588.4 | < 1s | < 1 |
| 10 | 10 | Hamming | 5725.4 | < 1s | < 1 |
| 10 | 10 | LeftM | 5725.4 | < 1s | < 1 |
| 10 | 100 | Hamming | 56517.4 | < 1s | < 1 |
| 10 | 100 | LeftM | 56517.4 | < 1s | < 1 |
| 10 | 250 | Hamming | 141851 | < 1s | 1 |
| 10 | 250 | LeftM | 141851 | < 1s | 1 |
| 100 | 1 | Hamming | 1035.2 | < 1s | < 1 |
| 100 | 1 | LeftM | 1035.2 | < 1s | < 1 |
| 100 | 10 | Hamming | 10435.4 | < 1s | < 1 |
| 100 | 10 | LeftM | 10435.4 | < 1s | < 1 |
| 100 | 100 | Hamming | 103598 | < 1s | 3 |
| 100 | 100 | LeftM | 103598 | < 1s | 3 |
| 100 | 250 | Hamming | 259113 | < 1s | 8 |
| 100 | 250 | LeftM | 259113 | < 1s | 8 |
| 1000 | 1 | Hamming | 1539.8 | < 1s | < 1 |
| 1000 | 1 | LeftM | 1539.8 | < 1s | < 1 |

| 1000 | 10 | Hamming | 14950.2 | < 1s | 4 |
|---|---|---|---|---|---|
| 1000 | 10 | LeftM | 14950.2 | < 1s | 4 |
| 1000 | 100 | Hamming | 149768 | < 1s | 38 |
| 1000 | 100 | LeftM | 149768 | < 1s | 38 |
| 1000 | 250 | Hamming | 375008 | 00:00:01 | 96 |
| 1000 | 250 | LeftM | 375008 | < 1s | 96 |
| 10000 | 1 | Hamming | 1928.2 | 00:00:04 | 6 |
| 10000 | 1 | LeftM | 1928.2 | 00:00:02 | 6 |
| 10000 | 10 | Hamming | 19569.2 | 00:00:15 | 50 |
| 10000 | 10 | LeftM | 19569.2 | 00:00:02 | 50 |
| 10000 | 100 | Hamming | 195747 | 00:02:05 | 492 |
| 10000 | 100 | LeftM | 195747 | 00:00:03 | 492 |
| 10000 | 250 | Hamming | 489204 | 00:05:13 | 1228 |
| 10000 | 250 | LeftM | 489204 | 00:00:03 | 1228 |

Table 2: Tree Sampling Threshold

| Threshold | # Markers | Time | Size [MByte] | SeU |
|---|---|---|---|---|
| 0.1 | 491025 | 00:00:24 | 1228 | 10.1 |
| 0.2 | 491025 | 00:00:56 | 1228 | 4.4 |
| 0.3 | 491025 | 00:01:30 | 1228 | 2.7 |
| 0.4 | 491025 | 00:02:02 | 1228 | 2 |
| 0.5 | 491025 | 00:02:37 | 1228 | 1.6 |
| 0.6 | 491025 | 00:02:56 | 1228 | 1.4 |
| 0.7 | 491025 | 00:03:41 | 1228 | 1.1 |
| 0.8 | 491025 | 00:04:05 | 1228 | 1 |
| 0.9 | 491025 | 00:04:16 | 1228 | 1 |
| 1 | 491025 | 00:04:06 | 1228 | — |

Speedups are computed with respect to a threshold of 1.

Table 3: ARG Window Width

| n | L [Mbp] | W [kbp] | # Markers | # Recs. | Time | Size [MByte] | SeU | SaU | RU |
|---|---|---|---|---|---|---|---|---|---|
| 10 | 1 | 0 | 549.2 | 171 | < 1s | < 1 | 1.2 | 1 | 1.1 |
| 10 | 1 | 100 | 549.2 | 176 | < 1s | < 1 | 1.1 | 1 | 1 |
| 10 | 1 | ∞ | 549.2 | 183 | < 1s | < 1 | — | — | — |
| 10 | 10 | 0 | 5694 | 1824 | < 1s | 1 | 1.1 | 1 | 1 |
| 10 | 10 | 100 | 5694 | 1875 | < 1s | 1 | 1.1 | 1 | 1 |
| 10 | 10 | ∞ | 5694 | 1811 | < 1s | 1 | — | — | — |
| 10 | 100 | 0 | 56502 | 19601 | 00:00:04 | 145 | 1.1 | 1 | 1 |
| 10 | 100 | 100 | 56502 | 19569 | 00:00:04 | 145 | 1.2 | 1 | 1 |
| 10 | 100 | ∞ | 56502 | 19591 | 00:00:05 | 144 | — | — | — |
| 10 | 250 | 0 | 141770 | 49098 | 00:00:27 | 893 | 1 | 1 | 1 |
| 10 | 250 | 100 | 141770 | 49210 | 00:00:25 | 894 | 1 | 1 | 1 |
| 10 | 250 | ∞ | 141770 | 49005 | 00:00:26 | 891 | — | — | — |
| 100 | 1 | 0 | 1064 | 1389 | < 1s | < 1 | 1.3 | 1 | 1 |
| 100 | 1 | 100 | 1064 | 1501 | < 1s | < 1 | 1.2 | 1 | 1 |
| 100 | 1 | ∞ | 1064 | 1455 | < 1s | < 1 | — | — | — |
| 100 | 10 | 0 | 10470 | 15163 | < 1s | 23 | 1.8 | 1 | 1 |
| 100 | 10 | 100 | 10470 | 15108 | < 1s | 22 | 1.7 | 1 | 1 |
| 100 | 10 | ∞ | 10470 | 15054 | 00:00:01 | 22 | — | — | — |
| 100 | 100 | 0 | 103530 | 149880 | 00:01:04 | 1989 | 1.7 | 1 | 1 |
| 100 | 100 | 100 | 103530 | 149580 | 00:01:07 | 1987 | 1.7 | 1 | 1 |
| 100 | 100 | ∞ | 103530 | 148700 | 00:01:52 | 1977 | — | — | — |
| 100 | 250 | 0 | 259080 | 375660 | 00:07:29 | 12403 | 1.5 | 1 | 1 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100 | 250 | 100 | 259080 | 375940 | 00:09:04 | 12396 | 1.2 | 1 | 1 |
| 100 | 250 | $\infty$ | 259080 | 373880 | 00:10:52 | 12324 | — | — | — |
| 1000 | 1 | 0 | 1473.8 | 6511 | $< 1s$ | 3 | 1.3 | 1 | 0.9 |
| 1000 | 1 | 100 | 1473.8 | 6554 | $< 1s$ | 3 | 1.2 | 1 | 0.9 |
| 1000 | 1 | $\infty$ | 1473.8 | 6032 | $< 1s$ | 3 | — | — | — |
| 1000 | 10 | 0 | 15109 | 68418 | 00:00:12 | 149 | 2.2 | 0.9 | 0.9 |
| 1000 | 10 | 100 | 15109 | 67575 | 00:00:13 | 148 | 2 | 1 | 1 |
| 1000 | 10 | $\infty$ | 15109 | 64855 | 00:00:26 | 141 | — | — | — |
| 1000 | 100 | 0 | 149680 | 671380 | 00:08:49 | 12744 | 2 | 1 | 1 |
| 1000 | 100 | 100 | 149680 | 666860 | 00:09:06 | 12722 | 2 | 1 | 1 |
| 1000 | 100 | $\infty$ | 149680 | 658390 | 00:17:47 | 12509 | — | — | — |
| 1000 | 250 | 0 | 374300 | 1658000 | 00:36:19 | 78668 | 1.7 | 1 | 1 |
| 1000 | 250 | 100 | 374300 | 1647900 | 00:37:19 | 78466 | 1.7 | 1 | 1 |
| 1000 | 250 | $\infty$ | 374300 | 1624700 | 01:03:06 | 77288 | — | — | — |
| 10000 | 1 | 0 | 1941.4 | 21168 | 00:00:27 | 16 | 1.1 | 1 | 0.9 |
| 10000 | 1 | 100 | 1941.4 | 21062 | 00:00:30 | 16 | 1 | 1 | 1 |
| 10000 | 1 | $\infty$ | 1941.4 | 20029 | 00:00:30 | 16 | — | — | — |
| 10000 | 10 | 0 | 19805 | 306650 | 00:10:10 | 870 | 1.3 | 0.8 | 0.8 |
| 10000 | 10 | 100 | 19805 | 299990 | 00:10:53 | 848 | 1.3 | 0.9 | 0.8 |
| 10000 | 10 | $\infty$ | 19805 | 251440 | 00:13:38 | 735 | — | — | — |
| 10000 | 100 | 0 | 195590 | 2927700 | 02:20:32 | 73383 | 1.6 | 0.9 | 0.9 |
| 10000 | 100 | 100 | 195590 | 2920100 | 02:40:46 | 72977 | 1.4 | 0.9 | 0.9 |
| 10000 | 100 | $\infty$ | 195590 | 2608600 | 03:46:53 | 65568 | — | — | — |
| 10000 | 250 | 0 | 490510 | 7170300 | 08:03:42 | 446480 | 2 | 0.9 | 0.9 |
| 10000 | 250 | 100 | 490510 | 7151900 | 08:39:40 | 445380 | 1.8 | 0.9 | 0.9 |
| 10000 | 250 | $\infty$ | 490510 | 6651000 | 15:56:16 | 415590 | — | — | — |

Speedups are computed with respect to an infinite window.

## 5.2 Secretary Sampler

The following result gives the probability that a run of the secretary sampler is exact.

**Lemma 1.** *Let $K$ be the index of the sample returned by a run of the secretary sampler on a vector of size $m$ with threshold $t_0$, $t_m = \lfloor mt_0 \rfloor$, $p_i = p_{ab_i}$, and $K_{t_m} = \arg\max_{i=1}^{t_m}\{g_i + \log p_i\}$. Let $\widehat{K}$ be the index of the correct sample. Then,*

$$\Pr(K = \widehat{K}) = \frac{1}{\sum_{i=1}^m p_i} \left( \sum_{i=1}^{t_m} p_i + p_{K_{t_m}}^{\widehat{K}-t_m-1} \left( \sum_{i=t_m+1}^m p_i \right) \prod_{i=t_m+1}^{\widehat{K}-1} (p_{K_{t_m}} + p_i)^{-1} \right).$$

*Proof.* By the formula of total probability,

$$\Pr(K = \widehat{K}) = \Pr(K = \widehat{K} \mid \widehat{K} \le t_m) \Pr(\widehat{K} \le t_m) + \Pr(K = \widehat{K} \mid \widehat{K} > t_m) \Pr(\widehat{K} > t_m).$$

Since the algorithm visits every element with an index below $t_m$, $\Pr(K = \widehat{K} \mid \widehat{K} \le t_m) = 1$. Given that $\widehat{K}$ is distributed as a categorical random variable, we obtain

$$\Pr(\widehat{K} \le t_m) = \frac{\sum_{i=1}^{t_m} p_i}{\sum_{i=1}^m p_i} \qquad\qquad \Pr(\widehat{K} > t_m) = \frac{\sum_{i=t_m+1}^m p_i}{\sum_{i=1}^m p_i}.$$

It only remains to show

$$\Pr(K = \widehat{K} \mid \widehat{K} > t_m) = p_{K_{t_m}}^{\widehat{K}-t_m-1} \prod_{i=t_m+1}^{\widehat{K}-1} (p_{K_{t_m}} + p_i)^{-1}.$$

By pairwise independence of the sequence $g_1, \ldots, g_m$, we obtain, after grouping similar terms, the following expression:

$$\Pr(K = \widehat{K} \mid \widehat{K} > t_m) = \prod_{i=t_m+1}^{\widehat{K}-1} \Pr(g_i - g_{K_{t_m}} \le \log p_{K_{t_m}} - \log p_i)$$

26

Since $g_i$ and $g_{K_{t_m}}$ follow a standard Gumbel distribution, their difference is distributed as a standard Logistic distribution. In consequence, the probability of interest is nothing more than the CDF of this distribution. Routine simplifications complete the proof. □

# References

[1] M. Nordborg. *Coalescent theory*, chapter 25, pages 843–877. Wiley, August 2007.

[2] Alexander L. Lewanski, Michael C. Grundler, and Gideon S. Bradburd. The era of the arg: an introduction to ancestral recombination graphs and their significance in empirical evolutionary genomics. *PLOS Genetics*, 20(1):e1011110, January 2024.

[3] J. F. C. Kingman. The coalescent. *Stochastic Processes and their Applications*, 13(3):235–248, September 1982.

[4] J. F. C. Kingman. On the genealogy of large populations. *Journal of Applied Probability*, 19(A):27–43, 1982.

[5] Richard R. Hudson. Properties of a neutral allele model with intragenic recombination. *Theoretical Population Biology*, 23(2):183–201, April 1983.

[6] Gilean A. T. McVean and Niall J. Cardin. Approximating the Coalescent with Recombination. *Philosophical Transactions: Biological Sciences*, 360(1459):1387–1393, 2005.

[7] Carsten Wiuf and Jotun Hein. Recombination as a point process along sequences. *Theoretical Population Biology*, 55(3):248–259, June 1999.

[8] Gary K. Chen, Paul Marjoram, and Jeffrey D. Wall. Fast and flexible simulation of DNA sequence data. *Genome Research*, 19(1):136–142, November 2008.

[9] Laurent Excoffier and Matthieu Foll. fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics*, 27(9):1332–1334, March 2011.

[10] Ying Wang, Ying Zhou, Linfeng Li, Xian Chen, Yuting Liu, Zhi-Ming Ma, and Shuhua Xu. A new method for modeling coalescent processes with recombination. *BMC Bioinformatics*, 15(1):273, August 2014.

[11] Paul R. Staab, Sha Zhu, Dirk Metzler, and Gerton Lunter. scrm: efficiently simulating long sequences using the approximated coalescent with recombination. *Bioinformatics*, 31(10):1680–1682, January 2015.

[12] Gregory Ewing and Joachim Hermisson. MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics*, 26(16):2064–2065, June 2010.

[13] Andrew D. Kern and Daniel R. Schrider. Discoal: flexible coalescent simulations with selection. *Bioinformatics*, 32(24):3839–3841, July 2016.

[14] Philipp W Messer. SLiM: Simulating Evolution with Selection and Linkage. *Genetics*, 194(4):1037–1039, August 2013.

[15] Jerome Kelleher, Alison M. Etheridge, and Gilean McVean. Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLOS Computational Biology*, 12(5):e1004842, May 2016.

[16] 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. Kretzschumar, Konrad Lohse, Michael Matschiner, Dominic Nelson, Nathaniel S. Pope, Consuelo D. Quinto-Cortés, Murillo F. Rodrigues, Kumar Saunack, Thibaut Sellinger, Kevin Thornton, Hugo van Kemenade, Anthony W. Wohns, 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), 12 2021.

[17] Gertjan Bisschop, Jerome Kelleher, and Peter L. Ralph. Likelihoods for a general class of args under the smc. *Genetics*, 2025.

[18] Fabrice Larribe, Sabin Lessard, and Nicholas J. Schork. Gene Mapping via the Ancestral Recombination Graph. *Theoretical Population Biology*, 62(2):215–229, September 2002.

[19] Patrick Fournier and Fabrice Larribe. Phenotype modelling using circuit theory on low treewidth networks. *To appear*, 2025.

[20] R. C. Griffiths and P. Marjoram. Ancestral inference from samples of DNA sequences with recombination. *Journal of Computational Biology: A Journal of Computational Molecular Cell Biology*, 3(4):479–502, 1996.

[21] Paul Fearnhead and Peter Donnelly. Estimating Recombination Rates From Population Genetic Data. *Genetics*, 159(3):1299–1318, November 2001.

[22] Jerome Kelleher, Yan Wong, Anthony W. Wohns, Chaimaa Fadil, Patrick K. Albers, and Gil McVean. Inferring whole-genome histories in large population datasets. *Nature Genetics*, 51(9):1330–1338, September 2019.

[23] Ali Mahmoudi, Jere Koskela, Jerome Kelleher, Yao-ban Chan, and David Balding. Bayesian inference of ancestral recombination graphs. *PLOS Computational Biology*, 18(3):e1009960, March 2022.

[24] Patrick Fournier and Fabrice Larribe. Human ancestries simulation and inference: a review of ancestral recombination graph samplers. *To appear*, 2025.

[25] Lusheng Wang, Kaizhong Zhang, and Louxin Zhang. Perfect phylogenetic networks with recombination. In *Proceedings of the 2001 ACM Symposium on Applied Computing*, pages 46–50, Las Vegas Nevada USA, March 2001. ACM.

[26] Jotun Hein. Reconstructing evolution of sequences subject to recombination using parsimony. *Mathematical Biosciences*, 98(2):185–200, March 1990.

[27] Yun S. Song and Jotun Hein. Constructing minimal ancestral recombination graphs. *Journal of Computational Biology*, 12(2):147–169, March 2005.

[28] Mark J. Minichiello and Richard Durbin. Mapping trait loci by use of inferred ancestral recombination graphs. *The American Journal of Human Genetics*, 79(5):910–922, November 2006.

[29] Thao Thi Phuong Nguyen, Vinh Sy Le, Hai Bich Ho, and Quang Si Le. Building ancestral recombination graphs for whole genomes. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 14(2):478–483, March 2017.

[30] Nguyen Thi Phuong Thao and Le Sy Vinh. A hybrid approach to optimize the number of recombinations in ancestral recombination graphs. In *Proceedings of the 2019 9th International Conference on Bioscience, Biochemistry and Bioinformatics*, ICBBB '19, pages 36–42. ACM, January 2019.

[31] Yan Wong, Anastasia Ignatieva, Jere Koskela, Gregor Gorjanc, Anthony W. Wohns, and Jerome Kelleher. A general and efficient representation of ancestral recombination graphs. *Genetics*, 228(1), July 2024.

[32] Emil Julius Gumbel. *Statistical Theory of Extreme Values and Some Practical Applications: A Series of Lectures*. U.S. Government Printing Office, 1954.

[33] Thomas S. Ferguson. Who Solved the Secretary Problem? *Statistical Science*, 4(3):282–289, August 1989.

[34] Erik Schnetter et al. *SIMD.jl*, 2025.

[35] Gaston H. Gonnet and Lawrence D. Rogers. The interpolation-sequential search algorithm. *Inf. Process. Lett.*, 6(4):136–139, 1977.

[36] Jeff Bonwick. The slab allocator: An object-caching kernel memory allocator. In *USENIX Summer*, pages 87–98. USENIX Association, 1994.

[37] Mason Protter et al. *Bumper.jl*, 2025.

[38] Benjamin Chu, Hua Zhou, Seyoon Ko, and Jcpapp. *VCFTools.jl*, 2023.

[39] Rune B. Lyngsø, Yun S. Song, and Jotun Hein. Minimum recombination histories by branch and bound. In *Algorithms in Bioinformatics*, pages 239–250. Springer Berlin Heidelberg, 2005.

[40] Anastasia Ignatieva, Rune B. Lyngsø, Paul A. Jenkins, and Jotun Hein. Kwarg: parsimonious reconstruction of ancestral recombination graphs with recurrent mutation. *Bioinformatics*, 37(19):3277–3284, May 2021.

[41] Fons van der Plas. *Pluto.jl*, 2025.

[42] Simon Danisch and Julius Krumbiegel. Makie.jl: flexible high-performance data visualization for julia. *Journal of Open Source Software*, 6(65):3349, September 2021.