

Efficient forwards-time simulation with ancestral recombination graphs

Jerome Kelleher, Kevin R. Thornton, Jaime Ashander, and Peter L. Ralph

October 23, 2017

Note: author order not determined

Abstract

To use genomic data for inference and prediction it is often necessary to obtain whole-genome information from individual-based simulations, but the computational burden of tracking the genome of each simulated individual can be substantial. In this note we describe how to both (a) dramatically reduce this burden and (b) efficiently record the entire history of the population. We do this by simulating only those loci that may affect reproduction (those having non-neutral variants), and recording the entire history of genetic inheritance in an efficient representation of the ancestral recombination graph, on which neutral mutations can be quickly placed afterwards. *make more clear data structure was already developed? refer to 'tree sequence' by name?* The algorithm is implemented in python, and is designed to be easily used by any forwards-time simulation software.

Coalescent simulations are very helpful but require random mating and neutrality. For continuous space, polygenic selection, or detailed dissection of life history, we must use forwards-time, individual-based simulation. These are much slower due in part to carrying around neutral genotypes irrelevant to the process. Here we show how to efficiently produce and store the entire history of ancestry and recombination (the ARG) from an individual-based simulation, on which neutral mutations can be placed afterwards. This has the promise of making large-scale, whole-genome simulations with realistic geography and selection finally possible.

OUTLINE

1. motivate need for whole-genome fwd-time simulations; point out that we only recently have the computing power to do this
2. explain ARG: explain that for forwards-time only need selected loci as by defn all others can be put on afterwards
3. review something about msprime methods for storing/traversing tree sequence
4. describe tables and write out conditions to have valid tables
5. write down algorithm used to do simple WF simulation
6. describe simplify algorithm
7. back-of-the-envelope calculation to compare cost of tracking whole genomes versus putting mutations on ARG
8. comparison of speed with simupop, fwdpp

Introduction

perhaps some of this veers into methods rather than introduction

Since the 1980's, coalescent theory has enabled computer simulation of the results of population genetics models identical to that which would be produced by large, randomly mating populations over long periods of time without actually requiring simulation of so many generations or meioses. Coalescent theory thus had three transformative effects on population genetics: first, giving researchers better conceptual tools to describe *gene trees* and thus bringing within-population trees into better focus; second, producing analytical methods to estimate parameters of interest from genetic data (e.g. $\theta = 4N_e\mu$); and finally, providing a computationally feasible method to produce computer simulations of population genetics processes. However, these powerful advances came with substantial caveats: the backwards-in-time processes that are described by coalescent theory are only *Markovian*, and thus feasible to work with, thanks to the important assumptions of (a) random mating, and (b) neutrality. *Brief statement why this is. Also include stationarity?* Both assumptions can be side-stepped to a limited extent, and so coalescent methods are now commonly used to simulate the results of population dynamics of a collection of randomly mating populations exchanging migrants, having a small number of loci under selection. Mapping results of such models onto real species can be challenging, as these are often distributed across geographical space and may have large numbers of loci under various sorts of selection. Furthermore, the relationship between the life history of a species – fecundity and mortality schedules, allee effects, and demographic fluctuations – are all absorbed into a single compound parameter, the coalescence rate. These considerations, and increasing computational power, have led to a resurgence of interest in forwards-time, individual-based simulations.

Another critical assumption of the coalescent is that the sample size n is much smaller than the effective population size N_e . With the very large sample sizes becoming common in human genetics, this is no longer a safe assumption. For example, a recent study [Martin et al., 2017] simulated 600,000 samples of human chromosome 20 to assess the impact of European biased reference panels in genome wide association studies. While many results of the coalescent are surprisingly robust to violations of the assumption that $n \ll N$, there are instances where genealogical properties are distorted and coalescent simulations may be misleading [Wakeley and Takahashi, 2003, Maruvka et al., 2011, Bhaskar et al., 2014]. A hybrid approach, in which the recent history of a large sample of individuals is simulated under a detailed forwards-in-time model and the deep history under the coalescent, is an attractive prospect. *Not sure where this should go; perhaps in the discussion? We can make the point that by keeping track of the genealogies forwards in time we can easily combine these with the deep history simulated by msprime*

With modern computing power, pure demographic calculations are not a barrier, even though biological population sizes are often above 10^6 , and coalescent theory tells us that a population of size N must be run for a multiple of N generations to produce stable genetic patterns. However, if our interest lies in the resulting genetic patterns of variation – and often, the point of such simulations is to compare to real data – then such simulations must somehow produce at the end data for each individual on a genomic scale. As samples of most species genomes harbor tens or hundreds of millions of variant sites, naively carrying full genotypes for even modest numbers of individuals through a simulation becomes quickly prohibitive.

However, it is thought that much of that variation is selectively neutral (or nearly so). By definition, the alleles carried by individuals in a population at neutral sites do not affect the population process. For this reason, if one records the entire genealogical history of a population over the course of a simulation, one can lay down neutral mutations on top of that history afterwards, without loss of generality. Precisely, we would need to know the genealogical tree relating all sampled individuals at each position along the genome. In this paper, we show how to use algorithmic tools and data structures developed for the coalescent simulator *msprime* to efficiently record, and later process, this history.

In so doing we record the *population pedigree* – the entire history of parent-offspring relationships of an entire population going back to a remote time – as well as information encoding the genetic outcomes of each ancestral meiosis – who inherited which parts of which parental chromosomes. This embellished pedigree contains all the information necessary to construct the genealogical tree that relates each individual to each other at each position on the genome, i.e., the *tree sequence*. Combined with ancestral genotypes and the origins of new mutations, it also completely specifies the genomic sequence of any individual in the

Figure 1: An example tree sequence with three samples over a chromosome of length 10. In the left-most panel we show the tree sequence pictorially in two different ways: in the top are the tree topologies and the bottom shows the spatial extent of the edges that define these topologies. The right-hand panels show the specific encoding of this tree sequence in the four tables (nodes, edges, sites and mutations) defined by `msprime`.

population at any time. This is much more than we need to know, however, so we discard all information irrelevant to the genetic history of the *sampled* individuals, which results in considerable savings.

Another way of representing this same information is known as the *ancestral recombination graph*, or ARG [Griffiths and Marjoram, 1997], which has been the subject of substantial study under the assumptions of coalescent theory [Wiuf and Hein, 1997, 1999, Marjoram and Wall, 2006, Wilton et al., 2015].

Methods

Reminder of what we need to know in the end (the trees), and quick review of `msprime` methods: sparse trees, tree differences.

Tree sequences

A tree sequence is an efficient encoding for a sequence of correlated trees such as is produced by recombination. The encoding is efficient because branches that are shared by adjacent trees are stored once, rather than repeatedly for each tree. The topology of a tree sequence is defined via a set of *nodes* and *edges*. A node in a tree sequence refers to a distinct ancestor and corresponds to the vertices in the individual genealogies along the sequence. Since each node represents a specific ancestor, it has a unique “time”, thought of as her birth time, which determines the height of any vertices she is associated with. In the example of Figure 1 we have a total of five nodes. Nodes 0, 1 and 2 occur at time 0 and are our samples. Nodes 3 and 4 represent the ancestors of these samples, and were born at time 1.5 and 2.5 respectively.

The topology of a tree sequence is defined by the *edges* which define how nodes relate to each other over specific genomic intervals. Each edge is a tuple (ℓ, r, p, c) , where $[\ell, r)$ is a half-open genomic interval defining the spatial extent of the edge; p and c are the parent and child nodes specifying a single branch in the tree over this interval. The spatial extent of the edges defining the topology of Figure 1 in the bottom left panel. For example, the branch joining nodes 1 to 3 is shared in both trees, and is recorded in a single edge extending over the whole chromosome. It is this method of capturing the shared structure between adjacent trees that makes the tree sequence encoding very compact and algorithmically efficient.

Continue along these lines with sites and mutations after updating the figure

Other stuff we can review here:

1. Encoding a tree sequence topology like this is very efficient, much, much more compact than newick. Can refer to example in `msprime` paper, or could do a large running example here. Say 500K samples of human CHR1, something like that. Newick will be hundreds of TB.
2. Can also make the point that it’s much more efficient than, say, VCF. We could drop some mutations on the example above, and then report how big it would be if encoded with VCF and BCF.
3. The `msprime` API is very efficient for working with these large datasets. Examples on how long it takes to iterate over all the trees, compute pi, etc.
4. The point of all this is to emphasise that not only is it efficient to generate the data as a tree sequence, it’s also very efficient to store and process it this way too. By the time we get to the end of this section we should have convinced the reader of this.

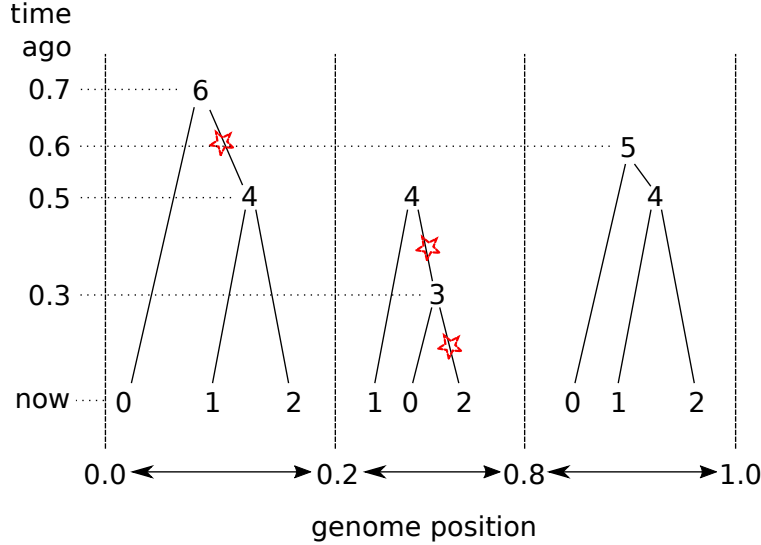


Figure 2: A pictorial representation of the ARG relating three samples to each other over a chromosome of length 1.0. Mutations shown in the example tables of the text are marked with red stars.

Data structures:

First we describe the data structures we use for recording genealogical history, as implemented in `msprime`. These derive from those described by Kelleher et al. [2016], but have been generalised and modified to remove redundancy. The tables below give the example tree sequence of Figure 2.

To clarify terminology, below a *tree* refers to a genealogical tree describing how a collection of individuals are related to each other. A *tree sequence* contains information sufficient to reconstruct the genealogical tree relating all samples to each other at any point along the genome. In the context of a tree sequence, *nodes* refer to distinct ancestors, and correspond to the vertices in the trees of a tree sequence. Since each node represents a certain ancestor, it has a unique “time”, thought of as her birth time, which determines the height of any branching points she is associated with. A given node will be associated with branching points of all trees across a region if that node is the most recent common ancestor to the subtending tips across that region. This information is stored in the columns of a **Node Table**:

id	is_sample	population	time
0	1	0	0
1	1	0	0
2	1	0	0
3	0	0	0.3
4	0	0	0.5
5	0	0	0.6
6	0	0	0.7

where “flags” records other information (e.g., a binary mask of ‘1’ indicates the node is a sample). Importantly, the **node ID** of a node is given implicitly by the (zero-based) index of its corresponding row in the Node Table.

Tree sequences are constructed by specifying over which segments of genome which nodes inherit from which other nodes. This information is stored by recording the endpoints of each distinctly inherited ancestral segment, the parental node, and a list of children nodes who have inherited that segment. As each such record describes a collection of edges across a swatch of trees in the tree sequence, we call these *edgesets* and store them in the columns of an **Edgeset Table**:

left	right	parent	children
0.2	0.8	3	0,2
0.0	0.2	4	1,2
0.2	0.8	4	1,3
0.8	1.0	4	1,2
0.8	1.0	5	0,4
0.0	0.2	6	0,4

To record information about genetic variants we need to also record each mutation and which nodes have inherited that mutation. The tree structure takes care of inheritance – all we need to do is to record the highest node in the tree at the mutated site that inherited that mutation. As more than one mutation may occur at a given site, we separate this information into two tables, first, the **Site Table** records for each variant site

id	position	ancestral_state
0	0.1	0
1	0.5	0

Here “position” is a (floating point) position along the chromosome, and “ancestral state” is the genotype of the root of the tree at that site. As for nodes, **site IDs** are given implicitly by the (zero-based) index of the rows. Then, we record in a **Mutation Table**

site	node	derived_state
0	4	1
1	3	1
1	2	0

in which “site” is the ID of the site at which this mutation occurred, “node” is the ID of the highest node that has inherited this mutation, and “derived state” is the genotype at this site of any individuals inheriting this mutation, unless another mutation occurs.

Definition of valid tables Here are the formal requirements for a set of nodes and edgesets to make sense, and to allow “msprime”’s algorithms to work properly.

To disallow time travel and multiple inheritance:

1. Offspring must be born after their parents (and hence, no loops).
2. The set of intervals on which each individual is a child must be disjoint.

For algorithmic reasons, we also require:

3. The leftmost endpoint of each chromosome is 0.0.
4. Node times must be strictly greater than zero.
5. The list of offspring in an edgeset must be sorted.
6. Edgesets must be sorted in nondecreasing time order.
7. The set of intervals on which each individual is a parent must be disjoint.
8. Each edgeset must contain at least two children.

Note that since each node time is equal to the amount of time since the *birth* of the corresponding parent, time is measured in clock time, not in meioses.

A forwards-time simulation does **not** naturally emit genealogical information satisfying requirements 5–8. However, **msprime** implements two algorithms that will take a set of tables satisfying only 1–4 and produce tables satisfying all requirements. Trivially, **sort_tables** enforces requirements 5 and 6 and does not renumber nodes; then, **simplify** enforces requirements 7 and 8 (and does much more; see below).

Recording the ARG in forwards time

To record the genealogical history of a simulation we need to record two things for each new chromosome: the birth time, and the endpoints and parental IDs of each distinctly inherited segment. This is recorded easily, without further processing, in the tables described above.

For concreteness, here we write out in pseudocode how to run a neutral Wright–Fisher simulation with overlapping generations that records genealogical history in this way. The simulation will run for T generations, and has N haploid individuals, each carrying a single chromosome of length L , on which for simplicity we assume there is exactly one crossover per generation. The probability of death per individual each generation is `death_prob`, and the whole-chromosome mutation rate per individual is `mut_rate`.

Initialize: We will build the tables `nodes`, `edgesets`, `sites`, and `mutations`, and keep track of the IDs of the current generation in the list `pop`.

```
for i in 0:N-1:
    nodes.add_row(time=T)
    pop[i] = i
```

(Alternatively, the tables could be initialized by the results of a coalescent simulation.)

Iterate: We then step through the generations, using the functions `random_mutation` and `random_allele` to choose positions and alleles for mutations, respectively:

```
for t in 1:T:
    for i in 0:N-1:
        if random.uniform() > death_prob:
            new_pop[i] = pop[i]
        else:
            u = nodes.num_rows # the ID of the new individual
            new_pop[i] = u
            nodes.add_row(time = T-t)
            a = random.sample(pop)
            b = random.sample(pop)
            bp = random.sample(0:L)
            if bp > 0:
                edgesets.add_row(left=0, right=bp,
                                parent=a, children=(u,))
            if bp < L:
                edgesets.add_row(left=bp, right=L,
                                parent=b, children=(u,))
            num_muts = random.poisson(mut_rate)
            for j in 1:num_muts:
                pos = random_mutation()
                if pos not in sites.position:
                    s = sites.num_rows
                    sites.add_row(position=pos,
                                  ancestral_state=random_allele())
                else:
                    s = sites.position.index(pos)
                    mutations.add_row(site=s, node=u,
                                      derived_state=random_allele())
    pop = new_pop
```

Finalize: to obtain a tree sequence, we need only transform the tables to the required format:

```
sort_tables(nodes, edgesets, sites, mutations)
simplify_tables(nodes, edgesets, sites, mutations)
```

which can then be loaded into a tree sequence with `msprime`. (These two functions operate on the tables in place.)

Tree sequence simplification

The resulting tables encode everything we need to know – in fact, they record all of history for everyone alive at any time through the simulation. This is much more than we need to reconstruct the genealogies and sequences of a smaller sample of individuals. Reducing this larger tree sequence to a smaller one relevant to a given set of “sample” nodes we call *simplification*. Roughly, this works by tracing ancestry from the samples backwards through the recorded history, adding node and edgeset records to the output only when coalescent events are reached. This works exactly as in `msprime`, allowing substantial re-use of algorithms; the main difference being that parental choice, mutations, and recombination locations are determined by the input tree sequence rather than randomly generated. An example is shown in Figure 3.

In this scheme, at any point in the simulation genealogical history is recorded in a tree sequence. This has two additional advantages. First, simplification can be run periodically through the simulation, taking the set of samples to be the entire currently alive population. This is important as it keeps memory usage from growing linearly (and quickly) with time. Second, the simulation can be *begun* with a tree sequence produced by some other method – for instance, by a coalescent simulation with `msprime`. This allows for incorporation of deep-time history beyond the reach of individual-based simulations. Since geographic structure from times longer ago than the mixing time of migration across the range has limited effect on modern genealogies [?] (other than possibly changing effective population size ?), this may not negatively affect realism.

Something about putting mutations down on the tree sequence.

Overview of the API

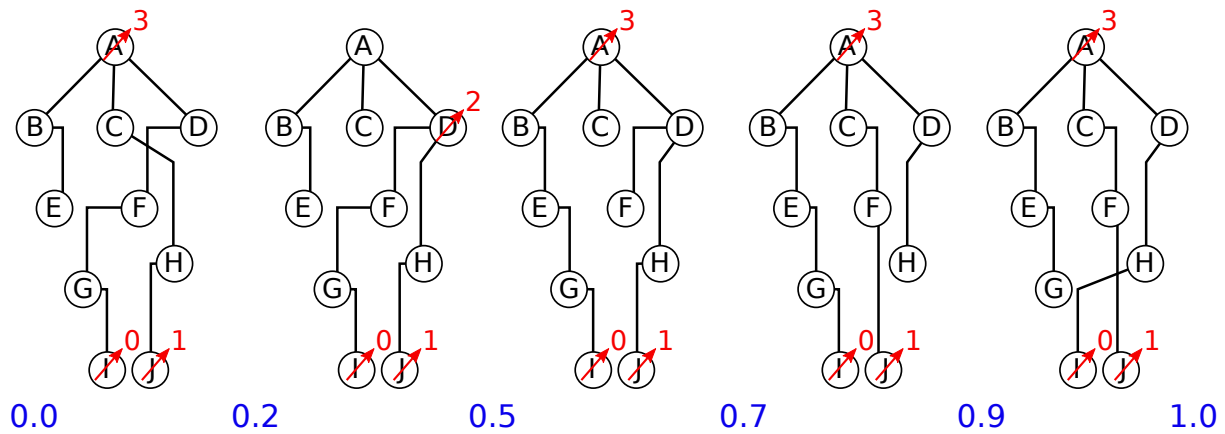
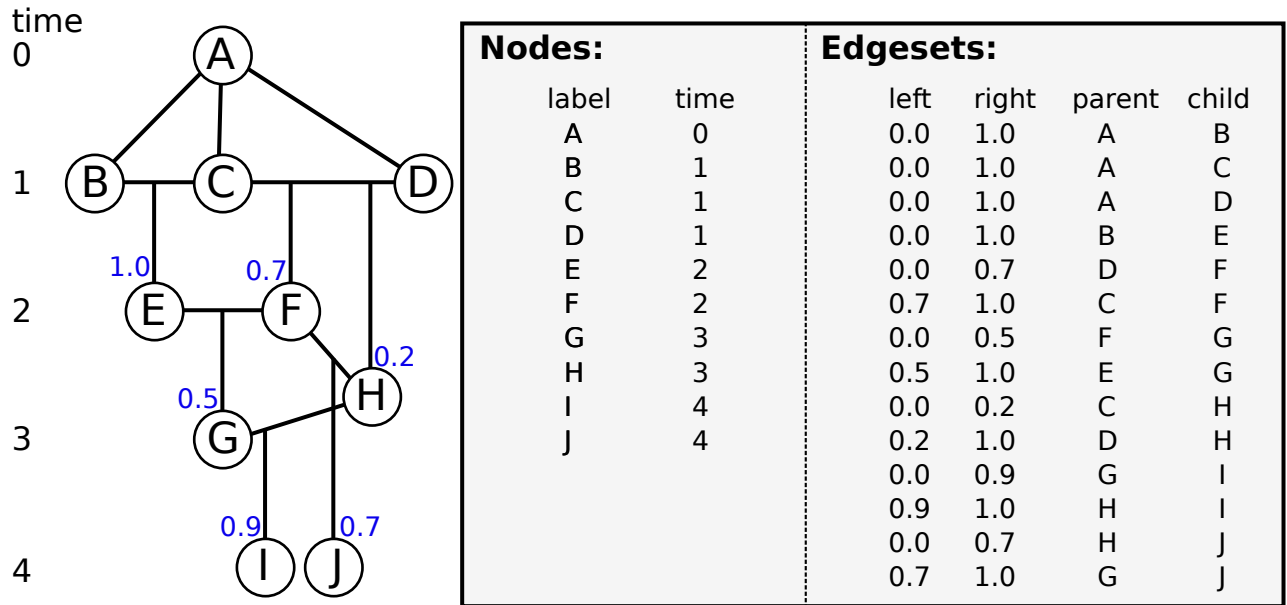
Quick overview of how to efficiently hook this up with other code.

Results

Estimates of run-time complexity Suppose that we wish to run a forwards-time simulation of N individuals for T generations, in which there are S selected loci and L neutral loci. We will estimate run-time complexity and memory usage for both a “naive” strategy that carries along neutral loci and an “ancestry-tracking” strategy like that we consider here. To do this, we assume that each individual must carry along its entire genotype. More advanced schemes are used in some simulators, but these increase efficiency by utilizing redundancy introduced by shared ancestry, which is effectively an intermediate scheme. We omit the cost of computing a fitness function.

Both schemes must choose mates and recombination breakpoints, and pass on selected genotypes. The difference between the two comes from the tradeoff between (a) passing on neutral genotypes, and (b) recording and simplifying the tree sequence, and adding neutral genotypes afterwards. (We assume here that selected genotypes are stored in the same way for both.) Passing on neutral genotypes naively records L items per individual each generation, discarding the previous generation. *do this better with numbers from nodes, edgesets* Recording a tree sequence stores 2 parents and ρ breakpoints on average each generation; after T generations this grows to $N \times T \times (2 + \rho)$ stored items. However, after simplification a tree sequence for all N individuals only takes of order $N + \rho \log N$ records. Simplification requires processing each of the initial records – so, order of $N \times T \times \rho$ operations;

Comparison of simulation with/without `msprime`, using `simuPOP` or maybe just a simple haploid simulation with 1000 QTL and stabilizing selection on a trait (say).



Simplified
for I,J:

Nodes:			Edgesets:			
	label	time ago	left	right	parent	children
I	0	0	0.0	0.2	3	0,1
J	1	0	0.2	0.5	2	0,1
D	2	3	0.5	1.0	3	0,1
A	3	4				

Figure 3: A simple example of the method. **Top:** the ARG shown on the left relates ten haploid individuals to each other. It is recorded, in forwards time, in 10 node records (one for each individual) and 14 edgeset records (one for each distinctly inherited segment). Blue numbers denote crossing over locations in each meiosis. The individuals *B*, *C*, and *D* inherit clonally from *A* to ensure rootedness of resulting tree. **Center:** the five distinct trees relating all individuals to each other found across the chromosome (blue numbers denote locations on the chromosome). Labels after simplification are shown in red. **Bottom:** tables recording the tree sequence after simplification with nodes *I* and *J* as samples. The mapping from labels in the forwards time simulation to nodes in the tree sequence is shown in red, which allows additional records to be added as the simulation progresses.

Maybe an estimate of how long *just* the ARG recording and simplification takes, so that then we can say how fast the simulator would have to be to do 10^6 whole chromosomes for 10^7 generations in a day.

Conclusion

This is a general-purpose strategy that can be applied to other methods.

All sorts of good reasons to want to have whole-genome simulations.

Acknowledgements

References

- A. Bhaskar, A. G. Clark, and Yun S. Song. Distortion of genealogical properties when the sample is very large. *Proc Natl Acad Sci U S A*, 111(6):2385–2390, 2014.
- Robert C. Griffiths and Paul Marjoram. An ancestral recombination graph. In *Progress in population genetics and human evolution (Minneapolis, MN, 1994)*, volume 87 of *IMA Vol. Math. Appl.*, pages 257–270. Springer, New York, 1997. URL <http://www.math.canterbury.ac.nz/~r.sainudiin/recomb/ima.pdf>.
- Jerome Kelleher, Alison M Etheridge, and Gilean McVean. Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS computational biology*, 12(5):e1004842, 2016.
- P Marjoram and J D Wall. Fast “coalescent” simulation. *BMC Genet*, 7:16–16, 2006. doi: 10.1186/1471-2156-7-16. URL <http://www.ncbi.nlm.nih.gov/pubmed/16539698>.
- Alicia R Martin, Christopher R Gignoux, Raymond K Walters, Genevieve L Wojcik, Benjamin M Neale, Simon Gravel, Mark J Daly, Carlos D Bustamante, and Eimear E Kenny. Human demographic history impacts genetic risk prediction across diverse populations. *The American Journal of Human Genetics*, 100(4):635–649, 2017.
- Yosef E Maruvka, Nadav M Shnerb, Yaneer Bar-Yam, and John Wakeley. Recovering population parameters from a single gene genealogy: an unbiased estimator of the growth rate. *Mol Biol Evol*, 28(5):1617–1631, 2011.
- John Wakeley and Tsuyoshi Takahashi. Gene genealogies when the sample size exceeds the effective size of the population. *Mol Biol Evol*, 20(2):208–213, 2003.
- P R Wilton, S Carmi, and A Hobolth. The SMC’ is a highly accurate approximation to the ancestral recombination graph. *Genetics*, 200(1):343–355, May 2015. doi: 10.1534/genetics.114.173898. URL <http://www.ncbi.nlm.nih.gov/pubmed/25786855>.
- C Wiuf and J Hein. On the number of ancestors to a DNA sequence. *Genetics*, 147(3):1459–1468, November 1997. URL <http://www.ncbi.nlm.nih.gov/pubmed/9383085>.
- C Wiuf and J Hein. The ancestry of a sample of sequences subject to recombination. *Genetics*, 151(3):1217–1228, March 1999. URL <http://www.ncbi.nlm.nih.gov/pubmed/10049937>.

A More general method for recording the ARG

moved from above

Concretely, this is done as follows. Suppose that the forwards simulation algorithm labels (haploid) individuals by integers, which we call “input labels”, to distinguish them from the “node IDs” given to these same individuals in the (output) tree sequence. The algorithm maintains at all times a set of tables (**nodes**, **edgesets**, **sites**, **mutations**) that record a relaxed tree sequence, and an associative array L that maps input labels to output node IDs, so that if x is an input label, then $L[x]$ is the corresponding output node ID. We also always maintain n to be the number of rows currently in the node table, (so that with zero-indexed IDs, the next to be added will have node ID n), T_0 to be the time of last simplification, and n_0 the number of rows in the node table at that time.

Initially, we begin with n and n_0 equal to the number of rows in the initial node table, and $L[j] = i_j$ for each $0 \leq j < N$ if the initial input generation is labeled i_0, \dots, i_{N-1} .

At a reproduction event where haploid parents x and y produce offspring u at time t in population p ,

1. add a (flags = 0, population = p , time = t) row to the node table,
2. set $L[u] = n$,
3. and increment $n += 1$.

Then, for each interval $[\ell, r)$ that u inherits from parent z (where z is either x or y),

4. add a (left = ℓ , right = r , parent = z , children = (u ,)) row to the edgeset table.

If furthermore there have been mutations at genomic locations s_1, \dots, s_k on this interval, with derived states q_1, \dots, q_k , then for each $1 \leq j \leq k$,

5. if s_j is not in the site table, add a row (position = s_j , ancestral_state = 0),
6. find the site index i whose position is s_j ,
7. and add a row (site = i , node= u , derived_state = q_j) to the mutation table.

To **simplify** the tree sequence at time T ,

1. add $T - T_0$ to each time in the first n_0 rows of the node table, and replace each remaining time t with $T - t$.

Then, pass the set of currently alive input IDs, i_0, \dots, i_{N-1} to the simplification algorithm, which produces a tree sequence that has node ID j corresponding to input ID i_j for $0 \leq j < N$, and

2. empty L ,
3. let $L[j] = i_j$ for $0 \leq j < N$,
4. set n to be the number of nodes in the new tree sequence and set $n_0 = n$, and finally
5. set $T_0 = T$.

Simplification keeps the tables to a manageable size. Since the map L is updated to maintain the association between individuals in the simulation and nodes in the tree sequence, simplification can be run regularly, as the simulation progresses.