



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Efficient ancestry and mutation simulation with msprime 1.0

**Citation for published version:**

Baumdicker, F, Bisschop, G, Goldstein, D, Gower, G, Ragsdale, AP, Tsambos, G, Zhu, S, Eldon, B, Ellerman, CE, Galloway, JG, Gladstein, AL, Gorjanc, G, Guo, B, Jeffery, B, Kretzschmar, WW, Lohse, K, Matschiner, M, Nelson, D, Pope, NS, Quinto-Cortés, CD, Rodrigues, MF, Saunack, K, Sellinger, T, Thornton, K, van Kemenade, H, Wohns, AW, Wong, HY, Gravel, S, Kern, AD, Koskela, J, Ralph, PL & Kelleher, J 2021, 'Efficient ancestry and mutation simulation with msprime 1.0', *Genetics*.  
<https://doi.org/10.1093/genetics/iyab229>

**Digital Object Identifier (DOI):**

[10.1093/genetics/iyab229](https://doi.org/10.1093/genetics/iyab229)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**Published In:**

Genetics

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Efficient ancestry and mutation simulation with msprime 1.0

Franz Baumdicker<sup>1,\*</sup>, Gertjan Bisschop<sup>2,\*</sup>, Daniel Goldstein<sup>3,24,\*</sup>, Graham Gower<sup>4,\*</sup>, Aaron P. Ragsdale<sup>5,\*</sup>, Georgia Tsambos<sup>6,\*</sup>, Sha Zhu<sup>7,\*</sup>, Bjarki Eldon<sup>8</sup>, E. Castedo Ellerman<sup>9</sup>, Jared G. Galloway<sup>10,11</sup>, Ariella L. Gladstein<sup>12,13</sup>, Gregor Gorjanc<sup>14</sup>, Bing Guo<sup>15</sup>, Ben Jeffery<sup>7</sup>, Warren W. Kretzschmar<sup>16</sup>, Konrad Lohse<sup>2</sup>, Michael Matschiner<sup>17</sup>, Dominic Nelson<sup>18</sup>, Nathaniel S. Pope<sup>19</sup>, Consuelo D. Quinto-Cortés<sup>20</sup>, Murillo F. Rodrigues<sup>10</sup>, Kumar Saunack<sup>21</sup>, Thibaut Sellinger<sup>22</sup>, Kevin Thornton<sup>23</sup>, Hugo van Kemenade<sup>24</sup>, Anthony W. Wohns<sup>7,25</sup>, Yan Wong<sup>7</sup>, Simon Gravel<sup>18,†</sup>, Andrew D. Kern<sup>10,†</sup>, Jere Koskela<sup>26,†</sup>, Peter L. Ralph<sup>10,27,†</sup> and Jerome Kelleher<sup>7,‡</sup>

<sup>1</sup>Cluster of Excellence “Controlling Microbes to Fight Infections”, Mathematical and Computational Population Genetics, University of Tübingen, 72076 Tübingen, Germany, <sup>2</sup>Institute of Evolutionary Biology, The University of Edinburgh, EH9 3FL, UK, <sup>3</sup>Khoury College of Computer Sciences, Northeastern University, MA 02115, USA, <sup>4</sup>Lundbeck GeoGenetics Centre, Globe Institute, University of Copenhagen, 1350 Copenhagen K, Denmark, <sup>5</sup>Department of Integrative Biology, University of Wisconsin–Madison, WI 53706, USA, <sup>6</sup>Melbourne Integrative Genomics, School of Mathematics and Statistics, University of Melbourne, Victoria, 3010, Australia, <sup>7</sup>Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford, OX3 7LF, UK, <sup>8</sup>Leibniz Institute for Evolution and Biodiversity Science, Museum für Naturkunde Berlin, 10115, Germany, <sup>9</sup>Fresh Pond Research Institute, Cambridge, MA 02140, USA, <sup>10</sup>Institute of Ecology and Evolution, Department of Biology, University of Oregon, OR 97403-5289, USA, <sup>11</sup>Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA 98102, USA, <sup>12</sup>Department of Genetics, University of North Carolina at Chapel Hill, NC 27599-7264, USA, <sup>13</sup>Embark Veterinary, Inc., Boston, MA 02111, USA, <sup>14</sup>The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of Edinburgh, EH25 9RG, UK, <sup>15</sup>Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD, 21201, USA, <sup>16</sup>Center for Hematology and Regenerative Medicine, Karolinska Institute, 141 83 Huddinge, Sweden, <sup>17</sup>Natural History Museum, University of Oslo, Blindern 0318 Oslo, Norway, <sup>18</sup>Department of Human Genetics, McGill University, Montréal, QC H3A 0C7, Canada, <sup>19</sup>Department of Entomology, Pennsylvania State University, PA 16802, USA, <sup>20</sup>National Laboratory of Genomics for Biodiversity (LANGEBIO), Unit of Advanced Genomics, CINVESTAV, Irapuato, Mexico, <sup>21</sup>IIT Bombay, Powai, Mumbai 400 076, Maharashtra, India, <sup>22</sup>Professorship for Population Genetics, Department of Life Science Systems, Technical University of Munich, 85354 Freising, Germany, <sup>23</sup>Ecology and Evolutionary Biology, University of California, Irvine, CA 92697, USA, <sup>24</sup>No affiliation, <sup>25</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA, <sup>26</sup>Department of Statistics, University of Warwick, CV4 7AL, UK, <sup>27</sup>Department of Mathematics, University of Oregon, OR 97403-5289 USA

**ABSTRACT** Stochastic simulation is a key tool in population genetics, since the models involved are often analytically intractable and simulation is usually the only way of obtaining ground-truth data to evaluate inferences. Because of this, a large number of specialized simulation programs have been developed, each filling a particular niche, but with largely overlapping functionality and a substantial duplication of effort. Here, we introduce msprime version 1.0, which efficiently implements ancestry and mutation simulations based on the succinct tree sequence data structure and the tskit library. We summarize msprime’s many features, and show that its performance is excellent, often many times faster and more memory efficient than specialized alternatives. These high-performance features have been thoroughly tested and validated, and built using a collaborative, open source development model, which reduces duplication of effort and promotes software quality via community engagement.

**KEYWORDS** Simulation, Coalescent, Mutations, Ancestral Recombination Graphs

## Introduction

The coalescent process (Kingman 1982a,b; Hudson 1983b; Tajima 1983) models the ancestry of a set of sampled genomes, providing a mathematical description of the genealogical tree that relates the samples to one another. It has proved to be a powerful model, and is now central to population genetics (Hudson 1990; Hein *et al.* 2004; Wakeley 2008). The coalescent is an ef-

doi: 10.1534/genetics.XXX.XXXXXX

Manuscript compiled: Wednesday 8<sup>th</sup> December, 2021

\*These authors contributed equally to this work.

†These authors contributed equally to this work.

‡jerome.kelleher@bdi.ox.ac.uk

© The Author(s) (2021). Published by Oxford University Press on behalf of the Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License

(<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium,

provided the original work is properly cited.

1 efficient framework for population genetic simulation, because  
 2 it allows us to simulate the genetic ancestry for a sample from  
 3 an idealized population model, without explicitly representing  
 4 the population in memory or stepping through the generations.  
 5 Indeed, [Hudson \(1983b\)](#) independently derived the coalescent  
 6 *in order to* efficiently simulate data, and used these simulations  
 7 to characterize an analytically intractable distribution. This in-  
 8 herent efficiency, and the great utility of simulations for a wide  
 9 range of purposes, has led to dozens of different tools being  
 10 developed over the decades ([Carvajal-Rodríguez 2008](#); [Liu et al.](#)  
 11 [2008](#); [Arenas 2012](#); [Yuan et al. 2012](#); [Hoban et al. 2012](#); [Yang et al.](#)  
 12 [2014](#); [Peng et al. 2015](#)).

13 Two technological developments of recent years, however,  
 14 pose major challenges to most existing simulation methods.  
 15 Firstly, fourth-generation sequencing technologies have made  
 16 complete chromosome-level assemblies possible ([Miga et al.](#)  
 17 [2020](#)), and high quality assemblies are now available for many  
 18 species. Thus, modeling genetic variation data as a series of  
 19 unlinked non-recombining loci is no longer a reasonable approx-  
 20 imation, and we must fully account for recombination. However,  
 21 while a genealogical tree relating  $n$  samples in the single-locus  
 22 coalescent can be simulated in  $O(n)$  time ([Hudson 1990](#)), the  
 23 coalescent with recombination is far more complex, and pro-  
 24 grams such as Hudson’s classical *ms* ([Hudson 2002](#)) can only  
 25 simulate short segments under the influence of recombination.  
 26 The second challenge facing simulation methods is that sam-  
 27 ple sizes in genetic studies have grown very quickly in recent  
 28 years, enabled by the precipitous fall in genome sequencing  
 29 costs. Human datasets like the UK Biobank ([Bycroft et al. 2018](#))  
 30 and gnomAD ([Karczewski et al. 2020](#)) now consist of hundreds  
 31 of thousands of genomes and many other datasets on a similar  
 32 scale are becoming available ([Tanjo et al. 2021](#)). Classical simu-  
 33 lators such as *ms* and even fast approximate methods such as  
 34 *scrm* ([Staab et al. 2015](#)) simply cannot cope with such a large  
 35 number of samples.

36 The *msprime* simulator ([Kelleher et al. 2016](#); [Kelleher and](#)  
 37 [Lohse 2020](#)) has greatly increased the scope of coalescent simu-  
 38 lations, and it is now straightforward to simulate millions of whole  
 39 chromosomes for a wide range of organisms. The “succinct tree  
 40 sequence” data structure ([Kelleher et al. 2016, 2018, 2019](#); [Wohns](#)  
 41 [et al. 2021](#)), originally introduced as part of *msprime*, makes it  
 42 possible to store such large simulations in a few gigabytes, sev-  
 43 eral orders of magnitude smaller than commonly used formats.  
 44 The succinct tree sequence has also led to major advances in  
 45 forwards-time simulation ([Kelleher et al. 2018](#); [Haller et al. 2018](#)),  
 46 ancestry inference ([Kelleher et al. 2019](#); [Wohns et al. 2021](#)) and  
 47 calculation of population genetic statistics ([Kelleher et al. 2016](#);  
 48 [Ralph et al. 2020](#)). Through a rigorous open-source community  
 49 development process, *msprime* has gained a large number of  
 50 features since its introduction, making it a highly efficient and  
 51 flexible platform for population genetic simulation. This paper  
 52 marks the release of *msprime* 1.0. We provide an overview of  
 53 its extensive features, demonstrate its performance advantages  
 54 over alternative software, and discuss opportunities for ongoing  
 55 open-source community-based development.

56 The efficiency of coalescent simulations depends crucially  
 57 on the assumption of neutrality, and it is important to note that  
 58 there are many situations in which this will be a poor approx-  
 59 imation of biological reality ([Johri et al. 2021](#)). In particular,  
 60 background selection has been shown to affect genome wide  
 61 sequence variation in a wide range of species ([Charlesworth et al.](#)  
 62 [1993, 1995](#); [Charlesworth and Jensen 2021](#)). Thus care must be

Interface	Separation of ancestry and mutation simulations. Ability to store arbitrary metadata along with simulation results, and automatic recording of provenance information for reproducibility. Jupyter notebook ( <a href="#">Kluyver et al. 2016</a> ) integration. Rich suite of analytical and visualization methods via the <i>tskit</i> library.
Ancestry	SMC, SMC’, Beta- and Dirac-coalescent, discrete time Wright-Fisher, and selective sweep models. Instantaneous bottlenecks. Discrete or continuous genomic coordinates, arbitrary ploidy, gene conversion. Output full ARG with recombination nodes, ARG likelihood calculations. Record full migration history and census events. Improved performance for large numbers of populations. Integration with forward simulators such as SLiM and fwdpy11 (“recapitation”).
Demography	Improved interface with integrated metadata and referencing populations by name. Import from Newick species tree, *BEAST ( <a href="#">Heled and Drummond 2009</a> ), and Demes ( <a href="#">Gower et al. 2022</a> ). Numerical methods to compute coalescence rates.
Mutations	JC69, HKY, F84, GTR, BLOSUM62, PAM, infinite alleles, SLiM and general matrix mutation models. Varying rates along the genome, recurrent/back mutations, discrete or continuous genomic coordinates, overlaying multiple layers of mutations, exact times associated with mutations.

**Table 1** Major features of *msprime* 1.0 added since version 0.3.0 ([Kelleher et al. 2016](#)).

63 taken to ensure that the results of purely neutral simulations are  
 64 appropriate for the question and genomic partition under study.  
 65 A major strength of *msprime*, however, is that it can be used in  
 66 conjunction with forwards-time simulators, enabling the simula-  
 67 tion of more realistic models than otherwise possible ([Kelleher](#)  
 68 [et al. 2018](#); [Haller et al. 2018](#)).

## 69 Results

70 In the following sections we describe the main features of  
 71 *msprime* 1.0, focusing on the aspects that are either new for  
 72 this version, or in which our approach differs significantly from  
 73 classical methods (summarized in Table 1). Where appropriate,  
 74 we benchmark *msprime* against other simulators, but the com-  
 75 parisons are illustrative and not intended to be systematic or  
 76 exhaustive. Please see [Kelleher et al. \(2016\)](#) for a performance  
 77 comparison of *msprime* against simulators such as *ms*, *msms*, and  
 78 *scrm*.

### 79 User interface

80 The majority of simulation packages are controlled either  
 81 through a command line interface (e.g. [Hudson 2002](#); [Kern and](#)  
 82 [Schridder 2016](#)), a text-based input file format (e.g. [Guillaume](#)  
 83 [and Rougemont 2006](#); [Excoffier and Foll 2011](#); [Shlyakhter et al.](#)  
 84 [2014](#)), or a mixture of both. Command line interfaces make it  
 85 easy to run simple simulations, but as model complexity and the  
 86 number of parameters increase, they become difficult to under-  
 87 stand and error-prone ([Ragsdale et al. 2020](#); [Gower et al. 2022](#)).  
 88 Specifying parameters through a text file alleviates this problem  
 89 to a degree, but lacks flexibility, for example, when running sim-  
 90 ulations with parameters drawn from a distribution. In practice,  
 91 for any reproducible simulation project users will write a script  
 92 to generate the required command lines or input parameter files,  
 93 invoke the simulation engine, and process the results in some

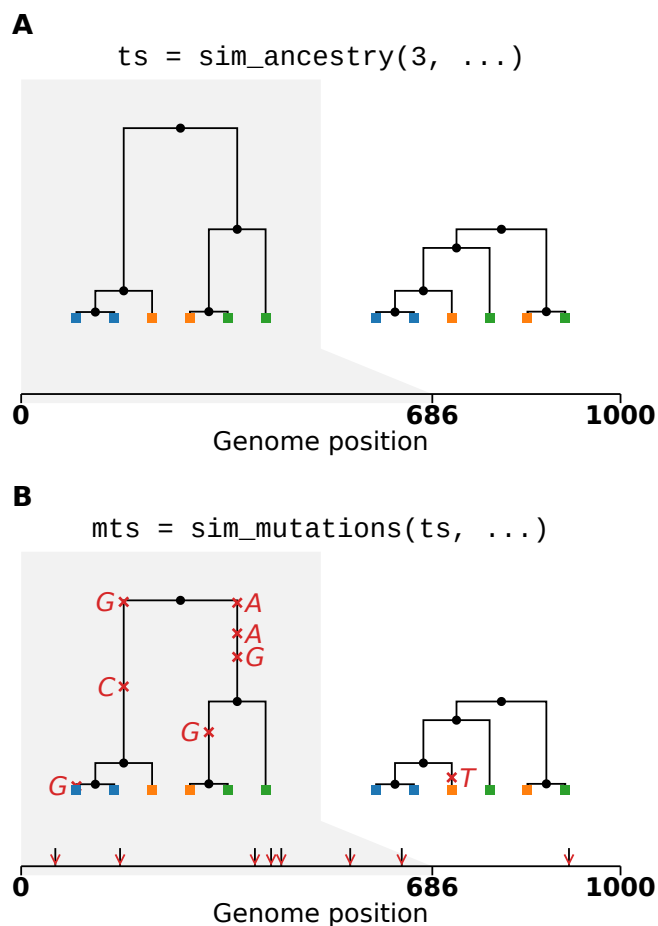
1 way. This process is cumbersome and labor intensive, and a  
 2 number of packages have been developed to allow simulations  
 3 to be run directly in a high-level scripting language (Staab and  
 4 Metzler 2016; Parobek *et al.* 2017; Gladstein *et al.* 2018).

5 The more recent trend has been to move away from this  
 6 file and command-line driven approach and to instead pro-  
 7 vide direct interfaces to the simulation engines via an Applica-  
 8 tion Programming Interface (API) (e.g. Thornton 2014; Kelleher  
 9 *et al.* 2016; Becheler *et al.* 2019; Haller and Messer 2019). The  
 10 primary interface for *msprime* is through a thoroughly docu-  
 11 mented Python API, which has encouraged the development of  
 12 an ecosystem of downstream tools (Terhorst *et al.* 2017; Chan *et al.*  
 13 2018; Spence and Song 2019; Adrion *et al.* 2020a,b; Kamm *et al.*  
 14 2020; McKenzie and Eaton 2020; Montinaro *et al.* 2020; Terasaki  
 15 Hart *et al.* 2021; Rivera-Colón *et al.* 2021). As well as providing  
 16 a stable and efficient platform for building downstream appli-  
 17 cations, *msprime*'s Python API makes it much easier to build  
 18 reproducible simulation pipelines, as the entire workflow can  
 19 be encapsulated in a single script, and package and version  
 20 dependencies explicitly stated using the pip or conda package  
 21 managers. For example, the errors made in the influential simu-  
 22 lation analysis of Martin *et al.* (2017) were only detected because  
 23 the pipeline could be easily run and reanalyzed (Ragsdale *et al.*  
 24 2020; Martin *et al.* 2020).

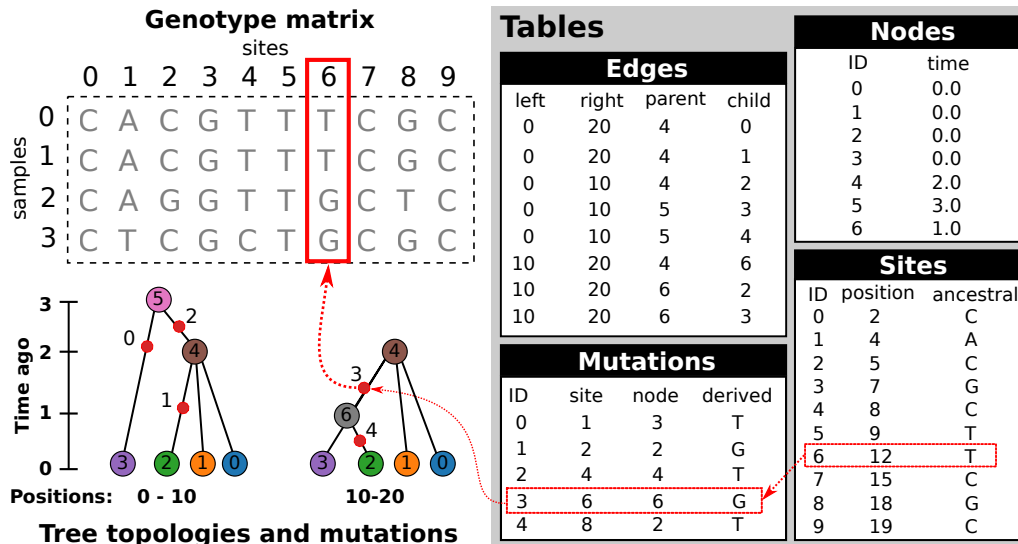
25 A major change for the *msprime* 1.0 release is the introduction  
 26 of a new set of APIs, designed in part to avoid sources of error  
 27 (see the Demography section) but also to provide more appro-  
 28 priate defaults while keeping compatibility with existing code.  
 29 In the new APIs, ancestry and mutation simulation are fully sep-  
 30 arated (see Fig. 1), with the `sim_ancestry` and `sim_mutations`  
 31 functions replacing the legacy `simulate` function. Among other  
 32 changes, the new APIs default to discrete genome coordinates  
 33 and finite sites mutations, making the default settings more real-  
 34 istic and resolving a major source of confusion and error. The  
 35 previous APIs are fully supported and tested, and will be main-  
 36 tained for the foreseeable future. The `msp` program (a command  
 37 line interface to the library) has been extended to include new  
 38 commands for simulating ancestry and mutations separately. A  
 39 particularly useful feature is the ability to specify demographic  
 40 models in Demes format (Gower *et al.* 2022) from the command  
 41 line, making simulation of complex demographies straightfor-  
 42 ward. We also provide an `ms` compatible command line interface  
 43 to support existing workflows.

#### 44 Tree sequences

45 One of the key reasons for *msprime*'s substantial performance  
 46 advantage over other simulators (Kelleher *et al.* 2016) is its use  
 47 of the "succinct tree sequence" data structure to represent simu-  
 48 lation results. The succinct tree sequence (usually abbreviated to  
 49 "tree sequence") was introduced by Kelleher *et al.* (2016) to con-  
 50 cisely encode genetic ancestry and sequence variation and was  
 51 originally implemented as part of *msprime*. We subsequently  
 52 extracted the core tree sequence functionality from *msprime* to  
 53 create the `tskit` library, which provides a large suite of tools for  
 54 processing genetic ancestry and variation data via APIs in the  
 55 Python and C languages (Tskit developers 2022). The availability  
 56 of `tskit` as a liberally licensed (MIT) open source toolkit has  
 57 enabled several other projects (e.g. Kelleher *et al.* 2019; Haller  
 58 and Messer 2019; Wohns *et al.* 2021; Terasaki Hart *et al.* 2021)  
 59 to take advantage of the same efficient data structures used in  
 60 *msprime*, and we hope that many more will follow. While a  
 61 full discussion of tree sequences and the capabilities of `tskit`



**Figure 1** Visualization of the separation between ancestry and mutation simulation. (A) The result of an invocation of `sim_ancestry` is two trees along a 1kb chunk of genome relating three diploid samples. Each diploid individual consists of two genomes (or nodes), indicated by color. (B) This ancestry is provided as the input to `sim_mutations`, which adds mutations. Graphics produced using `tskit`'s `draw_svg` method.



**Figure 2** An example tree sequence describing genealogies and sequence variation for four samples at ten sites on a chromosome of twenty bases long. Information is stored in a set of tables (the tables shown here include only essential columns, and much more information can be associated with the various entities). The node table stores information about sampled and ancestral genomes. The edge table describes how these genomes are related along a chromosome, and defines the genealogical tree at each position. The site and mutation tables together describe sequence variation among the samples. The genotype matrix and tree topologies shown on the left are derived from these tables.

1 is beyond the scope of this article, we summarize some aspects  
2 that are important for simulation.

3 Let us define a genome as the complete set of genetic material  
4 that a child inherits from one parent. Thus, a diploid individual  
5 has two (monoploid) genomes, one inherited from each parent.  
6 Since each diploid individual lies at the end of two distinct lin-  
7 eages of descent, they will be represented by *two* places (nodes)  
8 in any genealogical tree. In the tree sequence encoding a *node*  
9 therefore corresponds to a single genome, which is associated  
10 with its creation time (and other optional information), and  
11 recorded in a simple tabular format (Fig. 2). Genetic inheritance  
12 between genomes (nodes) is defined by edges. An *edge* consists  
13 of a parent node, a child node and the left and right coordinates  
14 of the contiguous chromosomal segment over which the child  
15 genome inherited genetic material from the parent genome. Par-  
16 ent and child nodes may correspond to ancestor and descendant  
17 genomes separated by many generations. Critically, edges can  
18 span multiple trees along the genome (usually referred to as  
19 “marginal” trees), and identical node IDs across different trees  
20 corresponds to the same ancestral genome. For example, in  
21 Fig. 2 the branch from node 0 to 4 is present in both marginal  
22 trees, and represented by a single edge (the first row in the edge  
23 table). This simple device, of explicitly associating tree nodes  
24 with specific ancestral genomes and recording the contiguous  
25 segments over which parent-child relationships exist, general-  
26 izes the original “coalescence records” concept (Kelleher *et al.*  
27 2016), and is the key to the efficiency of tree sequences (Kelleher  
28 *et al.* 2018, 2019; Ralph *et al.* 2020). Note that this formulation  
29 is fully compatible with the concept of an Ancestral Recombi-  
30 nation Graph (ARG) and any ARG topology can be fully and  
31 efficiently encoded in the node and edge tables illustrated in  
32 Fig. 2; see the section below for more details.

33 The final output of most population genetic simulations is  
34 some representation of sequence variation among the specified

35 samples. For coalescent simulations, we usually have three  
36 steps: (1) simulate the genetic ancestry, and optionally output  
37 the resulting marginal trees; (2) simulate sequence evolution  
38 conditioned on this ancestry by generating mutations (see the  
39 section); and (3) output the resulting nucleotide sequences by  
40 percolating the effects of the mutations through the trees. In-  
41 formation about the mutations themselves—e.g., where they  
42 have occurred on the trees—is usually not retained or made  
43 available for subsequent analysis. In *msprime*, however, we skip  
44 step (3), instead using *tskit*’s combined data model of ancestry  
45 and mutations to represent the simulated sequences. As illus-  
46 trated in Fig. 2, mutations are a fully integrated part of *tskit*’s  
47 tree sequence data model, and genetic variation is encoded by  
48 recording sites at which mutations have occurred, and where  
49 each mutation at those sites has occurred on the marginal tree.  
50 Crucially, the genome sequences themselves are never stored,  
51 or indeed directly represented in memory (although *tskit* can  
52 output the variant matrix in various formats, if required). It may  
53 at first seem inconvenient to have only this indirect represen-  
54 tation of the genome sequences, but it is extremely powerful.  
55 Firstly, the storage space required for simulations is dramatically  
56 reduced. For a simulation of  $n$  samples with  $m$  variant sites,  
57 we would require  $O(nm)$  space to store the sequence data as a  
58 variant matrix. However, if this simulation was of a recombining  
59 genome with  $t$  trees, then the *tskit* tree sequence encoding  
60 requires  $O(n + t + m)$  space, assuming we have  $O(1)$  mutations  
61 at each site (Kelleher *et al.* 2016). For large sample sizes, this  
62 difference is profound, making it conceivable, for example, to  
63 store the genetic ancestry and variation data for the entire hu-  
64 man population on a laptop (Kelleher *et al.* 2019). As well as  
65 the huge difference in storage efficiency, it is often far more ef-  
66 ficient to compute statistics of the sequence data from the trees  
67 and mutations than it is to work with the sequences themselves.  
68 For example, computing Tajima’s  $D$  from simulated data stored

1 in the `tskit` format is several orders of magnitude faster than  
2 efficient variant matrix libraries for large sample sizes (Ralph  
3 *et al.* 2020).

4 The vast genomic datasets produced during the SARS-CoV-2  
5 pandemic have highlighted the advantages of storing genetic  
6 variation data using the underlying trees. Turakhia *et al.* (2021)  
7 propose the Mutation Annotated Tree (MAT) format (consist-  
8 ing of a Newick tree and associated mutations in a binary for-  
9 mat) and the `matUtils` program as an efficient way to store and  
10 process large viral datasets (McBroome *et al.* 2021), achieving  
11 excellent compression and processing performance. Similarly,  
12 `phastsim` (De Maio *et al.* 2021) was developed to simulate se-  
13 quence evolution on such large SARS-CoV-2 phylogenies, and  
14 also outputs a Newick tree annotated with mutations (not in  
15 MAT format) to avoid the bottleneck of generating and storing  
16 the simulated sequences. While these methods illustrate the  
17 advantages of the general approach of storing ancestry and mu-  
18 tations rather than sequences, they do not generalize beyond  
19 their immediate settings, and no software library support is  
20 available.

21 The software ecosystem built around `tskit` is stable, mature  
22 and rapidly growing. Simulators such as `fdpdy11` (Thornton  
23 2014), `SLiM` (Haller and Messer 2019), `stdppsim` (Adrión *et al.*  
24 2020a), `Geonomics` (Terasaki Hart *et al.* 2021) and `GSpace` (Vir-  
25 goulay *et al.* 2021), and inference methods such as `tsinfer` (Kelle-  
26 her *et al.* 2019), `tsdate` (Wohns *et al.* 2021) and `Relate` (Speidel  
27 *et al.* 2019) use either the Python or C APIs to support outputting  
28 results in tree sequence format. Tree sequences are stored in an ef-  
29 ficient binary file format, and are fully portable across operating  
30 systems and processor architectures. The `tskit` library ensures  
31 interoperability between programs by having strict definitions  
32 of how the information in each of the tables is interpreted, and  
33 stringent checks for the internal consistency of the data model.

### 34 Data analysis

35 The standard way of representing simulation data is to render  
36 the results in a text format, which must subsequently be parsed  
37 and processed as part of some analysis pipeline. For example,  
38 `ms` outputs a set of sequences and can also optionally output the  
39 marginal trees along the genome in Newick format, and variants  
40 of this approach are used by many simulators. Text files have  
41 many advantages, but are slow to process at scale. The ability to  
42 efficiently process simulation results is particularly important  
43 in simulation-based inference methods such as Approximate  
44 Bayesian Computation (ABC) (Beaumont *et al.* 2002; Csilléry  
45 *et al.* 2010; Wegmann *et al.* 2010) and machine learning based  
46 approaches (Sheehan and Song 2016; Chan *et al.* 2018; Schrider  
47 and Kern 2018; Fligel *et al.* 2019; Sanchez *et al.* 2021). Clearly,  
48 simulation efficiency is crucial since the size and number of sim-  
49 ulations that can be performed determines the depth to which  
50 one can sample from the model and parameter space. Equally  
51 important, however, is the efficiency with which the simulation  
52 results can be transformed into the specific input required by  
53 the inference method. In the case of ABC, this is usually a set of  
54 summary statistics of the sequence data, and methods avoid the  
55 bottleneck of parsing text-based file formats to compute these  
56 statistics by either developing their own simulators (e.g. Cor-  
57 nuet *et al.* 2008; Lopes *et al.* 2009) or creating forked versions  
58 (i.e., modified copies) of existing simulators (e.g. Thornton and  
59 Andolfatto 2006; Hickerson *et al.* 2007; Pavlidis *et al.* 2010; Huang  
60 *et al.* 2011; Quinto-Cortés *et al.* 2018), tightly integrated with the  
61 inference method. Modern approaches to ABC such as ABC-

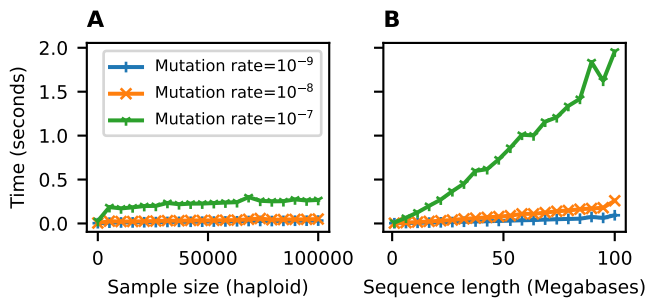
RF (Raynal *et al.* 2019; Pudlo *et al.* 2016) and ABC-NN (Csilléry  
*et al.* 2012; Blum and François 2010) use large numbers of weakly  
informative statistics, making the need to efficiently compute  
statistics from simulation results all the more acute. By using  
the stable APIs and efficient data interchange mechanisms pro-  
vided by `tskit`, the results of an `msprime` simulation can be  
immediately processed, without format conversion overhead.  
The `tskit` library has a rich suite of population genetic statistics  
and other utilities, and is in many cases orders of magnitude  
faster than matrix-based methods for large sample sizes (Ralph  
*et al.* 2020). Thus, the combination of `msprime` and `tskit` sub-  
stantially increases the overall efficiency of many simulation  
analysis pipelines.

Classical text based output formats like `ms` are inefficient to  
process, but also lack a great deal of important information about  
the simulated process. The tree-by-tree topology information  
output by simulators in Newick format lacks any concept of  
node identity, and means that we cannot reliably infer informa-  
tion about ancestors from the output. Because Newick stores  
branch lengths rather than node times, numerical precision is-  
sues also arise for large trees (McGill *et al.* 2013). Numerous  
forks of simulators have been created to access information  
not provided in the output. For example, `ms` has been forked  
to output information about migrating segments (Rosenzweig  
*et al.* 2016), ancestral lineages (Chen and Chen 2013), and `ms`'s  
fork `msHOT` (Hellenthal and Stephens 2007) has in turn been  
forked to output information on local ancestry (Racimo *et al.*  
2017). All of this information is either directly available by de-  
fault in `msprime`, or can be optionally stored via options such as  
`record_migrations` or `record_full_arg` (see the section) and  
can be efficiently and conveniently processed via `tskit` APIs.

### 52 Simulating mutations

53 Because coalescent simulations are usually concerned with neu-  
54 tral evolution (see the section, however) the problem of gener-  
55 ating synthetic genetic variation can be decomposed into two  
56 independent steps: firstly, simulating genetic ancestry (the trees),  
57 then subsequently simulating variation by superimposing muta-  
58 tion processes on those trees (see Fig. 1). A number of pro-  
59 grams exist to place mutations on trees: for instance, the classical  
60 `Seq-Gen` program (Rambaut and Grassly 1997) supports a range  
61 of different models of sequence evolution, and various exten-  
62 sions to the basic models have been proposed (e.g. Cartwright  
63 2005; Fletcher and Yang 2009). Partly for efficiency and partly in  
64 the interest of simplicity for users (i.e., to avoid intermediate text  
65 format conversions), population genetic simulators have tended  
66 to include their own implementations of mutation simulation,  
67 with most supporting the infinite sites model (e.g. Hudson 2002)  
68 but with several supporting a wide range of different models of  
69 sequence evolution (e.g. Mailund *et al.* 2005; Excoffier and Foll  
70 2011; Virgoulay *et al.* 2021). Thus, despite the logical separation  
71 between the tasks of simulating ancestry and neutral sequence  
72 evolution, the two have been conflated in practice.

Part of the reason for this poor record of software reuse and  
modularity is the lack of standardized file formats, and in par-  
ticular, the absence of common library infrastructure to abstract  
the details of interchanging simulation data. Although `msprime`  
also supports simulating both ancestry and mutations, the two  
aspects are functionally independent within the software; both  
ancestry and mutation simulators are present in `msprime` for rea-  
sons of convenience and history, and could be split into separate  
packages. The efficient C and Python interfaces for `tskit` make



**Figure 3** Time required to run `sim_mutations` on tree sequences generated by `sim_ancestry` (with a population size of  $10^4$  and recombination rate of  $10^{-8}$ ) for varying (haploid) sample size and sequence length. We ran 10 replicate mutation simulations each for three different mutation rates, and report the average CPU time required (Intel Core i7-9700). (A) Holding sequence length fixed at 10 megabases and varying the number of samples (tree tips) from 10 to 100,000. (B) Holding number of samples fixed at 1000, and varying the sequence length from 1 to 100 megabases. 250.14749pt

it straightforward to add further information to an existing file, and because of its efficient data interchange mechanisms, there is no performance penalty for operations being performed in a different software package. Thanks to this interoperability, `msprime`'s mutation generator can work with *any* `tskit` tree sequence, be it simulated using SLiM (Haller and Messer 2019) or `fdwp11` (Thornton 2014), or estimated from real data (Kelleher et al. 2019; Speidel et al. 2019; Wohns et al. 2021). It is a modular component intended to fit into a larger software ecosystem, and is in no way dependent on `msprime`'s ancestry simulator.

We have greatly extended the sophistication of `msprime`'s mutation generation engine for version 1.0, achieving near feature-parity with `Seq-Gen`. We support a large number of mutation models, including the JC69 (Jukes et al. 1969), F84 (Felsenstein and Churchill 1996), and GTR (Tavaré et al. 1986) nucleotide models and the BLOSUM62 (Henikoff and Henikoff 1992) and PAM (Dayhoff et al. 1978) amino acid models. Other models, such as the Kimura two and three parameter models (Kimura 1980, 1981), can be defined easily and efficiently in user code by specifying a transition matrix between any number of alleles. Mutation rates can vary along the genome, and multiple mutation models can be imposed on a tree sequence by overlaying mutations in multiple passes. We have extensively validated the results of mutation simulations against both theoretical expectations and output from `Seq-Gen` (Rambaut and Grassly 1997) and `Pyvolve` (Spielman and Wilke 2015).

Simulating mutations in `msprime` is efficient. Fig. 3 shows the time required to generate mutations (using the default JC69 model) on simulated tree sequences for a variety of mutation rates as we vary the number of samples (Fig. 3A) and the sequence length (Fig. 3B). For example, the longest running simulation in Fig. 3B required less than 2 seconds to generate an average of 1.5 million mutations over 137,081 trees in a tree sequence with 508,125 edges. This efficiency for large numbers of trees is possible because the tree sequence encoding allows us to generate mutations on an edge-by-edge basis (see Fig. 2 and the appendix), rather than tree-by-tree and branch-by-branch as would otherwise be required. Simulating mutations on a single tree is also very efficient; for example, we simulated mu-

tations under the BLOSUM62 amino acid model for a tree with  $10^6$  leaves over  $10^4$  sites (resulting in  $\sim 260,000$  mutations) in about 0.8 seconds, including the time required for file input and output. We do not attempt a systematic benchmarking of `msprime`'s mutation generation code against other methods, because at this scale it is difficult to disentangle the effects of inefficient input and output formats from the mutation generation algorithms. Given the above timings, it seems unlikely that generating mutations with `msprime` would be a bottleneck in any realistic analysis.

There are many ways in which the mutation generation code in `msprime` could be extended. For example, we intend to add support for microsatellites (Mailund et al. 2005), codon models (Arenas and Posada 2007) and indels (Cartwright 2005; Fletcher and Yang 2009), although changes may be required to `tskit`'s data model which is currently based on the assumption of independent sites.

### Recombination

Crossover recombination is implemented in `msprime` using Hudson's algorithm, which works backwards in time, generating common ancestor and recombination events and tracking their effects on segments of ancestral material inherited from the sample (Hudson 1983a, 1990; Kelleher et al. 2016). Common ancestor events merge the ancestral material of two lineages, and result in coalescences in the marginal trees when ancestral segments overlap. Recombination events split the ancestral material for some lineage at a breakpoint, creating two independent lineages. Using the appropriate data structures (Kelleher et al. 2016), this process is much more efficient to simulate than the equivalent left-to-right approach (Wiuf and Hein 1999b,a). In `msprime` 1.0, recombination rates can vary along a chromosome, allowing us to simulate recombination hotspots and patterns of recombination from empirical maps. The implementation of recombination in `msprime` is extensively validated against analytical results (Hudson 1983a; Kaplan and Hudson 1985) and simulations by `ms`, `msHOT` and SLiM.

The Sequentially Markovian Coalescent (SMC) is an approximation of the coalescent with recombination (McVean and Cardin 2005; Marjoram and Wall 2006), and was primarily motivated by the need to simulate longer genomes than was possible using tools like `ms`. The SMC is a good approximation to the coalescent with recombination when we have fewer than five sampled genomes (Hobolth and Jensen 2014; Wilton et al. 2015), but the effects of the approximation are less well understood for larger sample sizes, and several approaches have been proposed that allow simulations to more closely approximate the coalescent with recombination (Chen et al. 2009; Wang et al. 2014; Staab et al. 2015). The SMC and SMC' models are supported in `msprime` 1.0. However, they are currently implemented using a naive rejection sampling approach, and are somewhat slower to simulate than the exact coalescent with recombination. These models are therefore currently only appropriate for studying the SMC approximations themselves, although we intend to implement them more efficiently in future versions.

In human-like parameter regimes and for large sample sizes, `msprime`'s implementation of the exact coalescent with recombination comprehensively outperforms all other simulators, including those based on SMC approximations (Kelleher et al. 2016). However, it is important to note that although the implementation of Hudson's algorithm is very efficient, it is still quadratic in the population scaled recombination rate  $\rho = 4N_eL$ ,

where  $L$  is the length of the genome in units of recombination distance. This is because Hudson’s algorithm tracks recombinations not only in segments ancestral to the sample, but also between ancestral segments. As mentioned above, common ancestor events in which the ancestral material of two lineages is merged only result in coalescences in the marginal trees if their ancestral segments overlap. If there is no overlap, the merged segments represent an ancestral chromosome that is a genetic ancestor of the two lineages, but not the most recent common genetic ancestor at any location along the genome. When this happens, the merged lineage carries “trapped” genetic material that is not ancestral to any samples, but where recombinations can still occur (Wiuf and Hein 1999b). For large  $\rho$ , recombination events in trapped ancestral material will dominate, and so we can use this as a proxy for the overall number of events in Hudson’s algorithm. Hein *et al.* (2004, Eq. 5.10) gave

$$\rho(\rho + 1) \left( \sum_{i=1}^{n-1} \frac{1}{i} \right)^2 \quad (1)$$

1 as an upper bound on the number of recombination events  
 2 within trapped ancestral material for  $n$  samples. As discussed in  
 3 the appendix, the quadratic dependence of simulation running  
 4 time on  $\rho$  implied by (1) is well supported by observations, and  
 5 provides a useful means of predicting how long a particular  
 6 simulation might require.

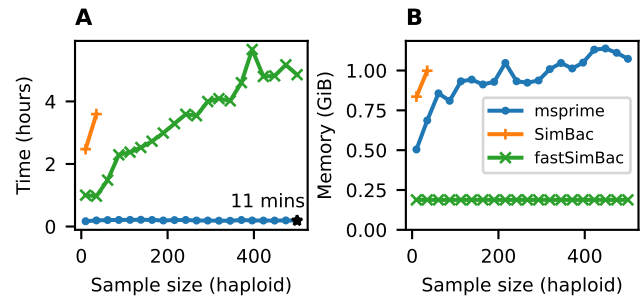
### 7 Gene conversion

8 Gene conversion is a form of recombination that results in the  
 9 transfer of a short segment of genetic material, for example be-  
 10 tween homologous chromosomes (Chen *et al.* 2007). Since gene  
 11 conversion impacts much shorter segments than crossover re-  
 12 combination (typically below 1kb) it affects patterns of linkage  
 13 disequilibrium differently (Korunes and Noor 2017). Wiuf and  
 14 Hein (2000) modeled gene conversion in the coalescent via a rate  
 15 at which gene conversion events are initiated along the genome  
 16 and a geometrically distributed tract length. In terms of the  
 17 ancestral process, gene conversion differs from crossover recom-  
 18 bination (as described in the previous section) in that it extracts  
 19 a short tract of ancestry into an independent lineage, rather than  
 20 splitting ancestry to the left and right of a given breakpoint. We  
 21 have implemented this model of gene conversion in *msprime* 1.0,  
 22 and validated the output against *ms* and analytical results (Wiuf  
 23 and Hein 2000).

24 Gene conversion is particularly useful to model homolo-  
 25 gous recombination in bacterial evolution, and so we compare  
 26 the performance of *msprime* with gene conversion to two spe-  
 27 cialized bacterial simulators, *SimBac* (Brown *et al.* 2016) and  
 28 *fastSimBac* (De Maio and Wilson 2017). Figure 4A shows  
 29 that *msprime* is far more efficient than both *SimBac* and the  
 30 SMC-based approximation *fastSimBac*. Figure 4B shows that  
 31 *msprime* requires somewhat more memory than *fastSimBac*,  
 32 (as expected since *fastSimBac* uses a left-to-right SMC approx-  
 33 imation) but is still reasonably modest at around 1GiB for a  
 34 simulation of 500 whole *E. coli* genomes. However, *msprime* is  
 35 currently lacking many of the specialized features required to  
 36 model bacteria, and so an important avenue for future work  
 37 is to add features such as circular genomes and bacterial gene  
 38 transfer (Baumdicker and Pfaffelhuber 2014).

### 39 Demography

40 One of the key applications of population genetic simulations is  
 41 to generate data for complex demographies. Beyond idealized



**Figure 4** Comparison of simulation performance using *msprime* (`sim_ancestry`), *SimBac*, and *fastSimBac* for varying (haploid) sample sizes, and the current estimates for *E. coli* parameters (Lapierre *et al.* 2016): a 4.6Mb genome,  $N_e = 1.8 \times 10^8$ , gene conversion rate of  $8.9 \times 10^{-11}$  per base and mean tract length of 542. We report (A) the total CPU time and (B) maximum memory usage averaged over 5 replicates (Intel Xeon E5-2680 CPU). We did not run *SimBac* beyond first two data points because of the very long running times.

cases such as stepping-stone or island models, or specialized cases such as isolation-with-migration models, analytical results are rarely possible. Simulation is therefore integral to the development and evaluation of methods for demographic inference. The demography model in *msprime* is directly derived from the approach used in *ms*, and supports an arbitrary number of randomly mating populations exchanging migrants at specified rates. A range of demographic events are supported, which allow for varying population sizes and growth rates, changing migration rates over time, as well as population splits, admixtures and pulse migrations.

A major change for *msprime* 1.0 is the introduction of the new Demography API, designed to address a design flaw in the *msprime* 0.x interface which led to avoidable errors in downstream simulations (Ragsdale *et al.* 2020). The new API is more user-friendly, providing the ability, for example, to refer to populations by name rather than their integer identifiers. We also provide numerical methods to compute the coalescence rates for two or more lineages which can be inverted to obtain the “inverse instantaneous coalescence rate” of Chikhi *et al.* (2018). Many popular approaches in population genetics use the distribution of coalescence rates between pairs of lineages to infer effective population sizes over time (Li and Durbin 2011; Sheehan *et al.* 2013; Schiffels and Durbin 2014) or split times and subsequent migration rates between populations (Wang *et al.* 2020). These numerical methods provide a valuable ground-truth when evaluating such inference methods, as illustrated by Adrion *et al.* (2020a).

### Instantaneous bottlenecks

A common approach to modeling the effect of demographic history on genealogies is to assume that effective population size ( $N_e$ ) changes in discrete steps which define a series of epochs (Griffiths *et al.* 1994; Marth *et al.* 2004; Keightley and Eyre-Walker 2007; Li and Durbin 2011). In this setting of piecewise constant  $N_e$ , capturing a population bottleneck requires three epochs:  $N_e$  is reduced by some fraction  $b$  at the start of the bottleneck,  $T_{start}$ , and recovers to its initial value at time  $T_{end}$  (Marth *et al.* 2004). If bottlenecks are short both on the timescale of coa-



1 lence and mutations, there may be little information about the  
 2 duration of a bottleneck ( $T_{end} - T_{start}$ ) in sequence data. Thus  
 3 a simpler, alternative model is to assume that bottlenecks are  
 4 instantaneous ( $T_{end} - T_{start} \rightarrow 0$ ) and generate a sudden burst  
 5 of coalescence events (a multiple merger event) in the genealogy.  
 6 The strength of the bottleneck  $B$  can be thought of as an  
 7 (imaginary) time period during which coalescence events are  
 8 collapsed, i.e. there is no growth in genealogical branches during  
 9  $B$  and the probability that a single pair of lineages entering the  
 10 bottleneck coalesce during the bottleneck is  $1 - e^{-B}$ . Although  
 11 this simple two parameter model of bottlenecks is attractive and  
 12 both analytic results and empirical inference (Griffiths *et al.*  
 13 1994; Birkner *et al.* 2009; Galtier *et al.* 2000; Bunnefeld *et al.* 2015)  
 14 have been developed under this model, there has been no software  
 15 available to simulate data under instantaneous bottleneck  
 16 histories.

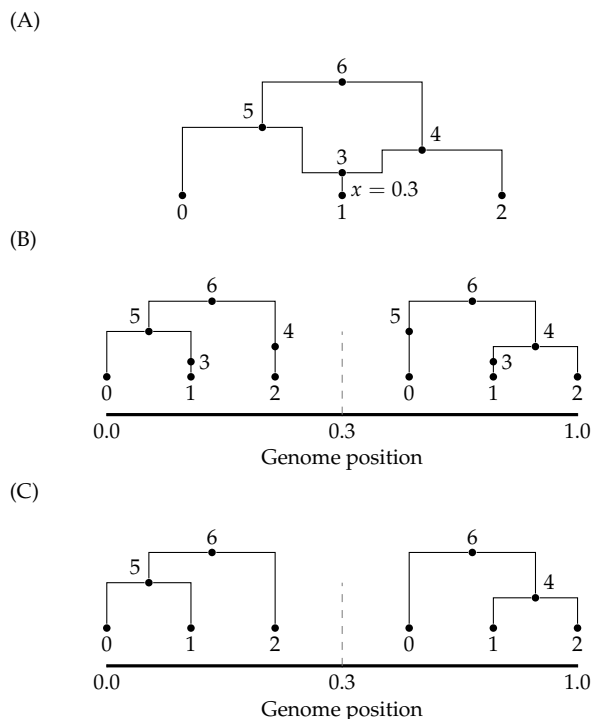
17 We have implemented instantaneous bottlenecks in  
 18 msprime 1.0 using a variant of Hudson’s linear time single-locus  
 19 coalescent algorithm (Hudson 1990), and validated the results  
 20 by comparing against analytical expectations (Bunnefeld *et al.*  
 21 2015).

### 22 Multiple merger coalescents

23 Kingman’s coalescent assumes that only two ancestral lineages  
 24 can merge at each merger event. Although this is generally a reasonable  
 25 approximation, there are certain situations in which the underlying  
 26 mathematical assumptions are violated. For example in certain highly  
 27 fecund organisms (Hedgecock 1994; Beckenbach 1994; Hedgecock and  
 28 Pudovkin 2011; Árnason 2004; Irwin *et al.* 2016; Vendrami *et al.* 2021),  
 29 where individuals have the ability to produce numbers of offspring on  
 30 the order of the population size and therefore a few individuals may  
 31 produce the bulk of the offspring in any given generation (Hedgecock  
 32 1994). These population dynamics violate basic assumptions of the  
 33 Kingman coalescent, and are better modeled by ‘multiple-merger’  
 34 coalescents (Donnelly and Kurtz 1999; Pitman 1999; Sagitov 1999;  
 35 Schweinsberg 2000; Möhle and Sagitov 2001), in which more than  
 36 two lineages can merge in a given event. Multiple-merger  
 37 coalescent processes have also been shown to be relevant for  
 38 modeling the effects of selection on gene genealogies (Gillespie  
 39 2000; Durrett and Schweinsberg 2004; Desai *et al.* 2013; Neher  
 40 and Hallatschek 2013; Schweinsberg 2017).

42 Although multiple merger coalescents have been of significant  
 43 theoretical interest for around two decades, there has been little  
 44 practical software available to simulate these models. Kelleher  
 45 *et al.* (2013, 2014) developed packages to simulate a related  
 46 spatial continuum model (Barton *et al.* 2010), Zhu *et al.* (2015)  
 47 simulate genealogies within a species tree based on a multiple-  
 48 merger model, and Becheler and Knowles (2020) provide a general  
 49 method for simulating multiple merger processes as part of the  
 50 Quetzal framework (Becheler *et al.* 2019). The Beta- $\Xi$ -Sim  
 51 simulator (Koskela 2018; Koskela and Wilke Berenguer 2019) also  
 52 includes a number of extensions to the Beta-coalescent. None of  
 53 these methods work with large genomes, and very little work  
 54 has been performed on simulating multiple merger processes  
 55 with recombination.

56 We have added two multiple merger coalescent models in  
 57 msprime 1.0, the Beta-coalescent (Schweinsberg 2003) and  
 58 “Dirac”-coalescent (Birkner *et al.* 2013a), allowing us to efficiently  
 59 simulate such models with recombination for the first time. These  
 60 simulation models have been extensively validated against analytical  
 61 results from the site frequency spec-



**Figure 5** (A) A simple ARG in which a recombination occurs at position 0.3; (B) the equivalent topology depicted as a tree sequence, including the recombination node; (C) the same tree sequence topology “simplified” down to its minimal tree sequence representation. Note the original node IDs have been retained for clarity.

trum (Birkner *et al.* 2013b; Blath *et al.* 2016; Hobolth *et al.* 2019) as well as more general properties of coalescent processes. See the appendix for more details and model derivations.

### 65 Ancestral Recombination Graphs

66 The Ancestral Recombination Graph (ARG) was introduced by Griffiths (Griffiths 1991; Griffiths and Marjoram 1997) to represent the stochastic process of the coalescent with recombination as a graph. This formulation is complementary to Hudson’s earlier work (Hudson 1983a), and substantially increased our theoretical understanding of recombination. In Griffiths’ ARG formulation, a realization of the coalescent with recombination is a graph in which vertices represent common ancestor or recombination events, and edges represent lineages. There is the “big” ARG, in which we track lineages arising out of recombinations regardless of whether they carry ancestral material (Ethier and Griffiths 1990), and the “little” ARG in which we only track genetic ancestors. Over time, usage of the term has shifted away from its original definition as a stochastic process, to being interpreted as a representation of a particular genetic ancestry as a graph, without necessarily following the specific details of the Griffiths formulation (e.g. Minichiello and Durbin 2006; Mathieson and Scally 2020). Under the latter interpretation, the tree sequence encoding of genetic ancestry (described above) clearly is an ARG: the nodes and edges define a graph in which edges are annotated with the set of disjoint genomic intervals through which ancestry flows.

68 For our purposes, an ARG is a realization of the coalescent with recombination, in the Griffiths (little ARG) sense. As de-

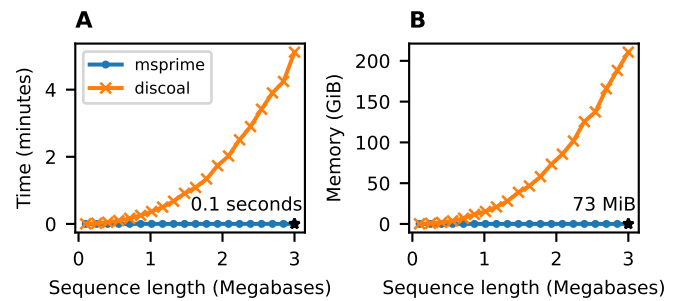
scribed in detail by Kelleher *et al.* (2016), Hudson’s algorithm works by dynamically traversing a little ARG. The graph is not explicitly represented in memory, but is partially present through the extant lineages and the ancestral material they carry over time. We do not output the graph directly, but rather store the information required to recover the genealogical history as nodes and edges in a tree sequence. This is far more efficient than outputting the simulated ARG in its entirety. For a given scaled recombination rate  $\rho$  (setting aside the dependency on the sample size  $n$ ) we know from Eq. (1) that the number of nodes in an ARG is  $O(\rho^2)$ , whereas the size of the tree sequence encoding is  $O(\rho)$  (Kelleher *et al.* 2016). This difference between a quadratic and a linear dependency on  $\rho$  is profound, and shows why large simulations cannot output an ARG in practice.

Although by default `msprime` outputs tree sequences that contain full information about the genealogical trees, their correlation structure along the chromosome, and the ancestral genomes on which coalescences occurred, some information is lost in this mapping down from ARG space to the minimal tree sequence form. In particular, we lose information about ancestral genomes that were common ancestors but in which no coalescences occurred, and also information about the precise time and chromosomal location of recombination events. In most cases, such information is of little relevance as it is in principle unknowable, but there are occasions such as visualization or computing likelihoods (see below) in which it is useful. We therefore provide the `record_full_arg` option in `msprime` to store a representation of the complete ARG traversed during simulation. This is done by storing extra nodes (marked with specific flags, so they can be easily identified later) and edges in the tree sequence (Fig. 5). One situation in which a record of the full ARG is necessary is when we wish to compute likelihoods during inference. The likelihood is a central quantity in evaluating the plausibility of a putative ancestry as an explanation of DNA sequence data, both directly through e.g. approaches based on maximum likelihood, and as an ingredient of methods such as Metropolis-Hastings (Kuhner *et al.* 2000; Nielsen 2000; Wang and Rannala 2008). We provide functions to compute the likelihood of ARG realizations and mutational patterns under the standard coalescent and infinite sites mutation model. For details, see the appendix: .

### Selective sweeps

Another elaboration of the standard neutral coalescent with recombination is the addition of selective sweeps (Kaplan *et al.* 1989; Braverman *et al.* 1995; Kim and Stephan 2002). Sweeps are modeled by creating a structured population during the sojourn of the beneficial mutation through the population (i.e., the sweep phase) in which lineages may transit between favored and unfavoured backgrounds through recombination. This approach allows for many selective sweep scenarios to be simulated efficiently, including recurrent, partial, and soft selective sweeps. However this efficiency comes at the cost of flexibility in comparison to forwards in time simulation. Several specialized simulators have been developed to simulate sweeps in the coalescent, including `SeLSim` (Spencer and Coop 2004), `mbs` (Teshima and Innan 2009), `msms` (Ewing and Hermisson 2010), `cosi2` (Shlyakhter *et al.* 2014) and `discoal` (Kern and Schrider 2016).

Selective sweeps are implemented in the coalescent as a two step-process: first generating an allele frequency trajectory, and then simulating a structured coalescent process conditioned on that trajectory. Following `discoal`, we generate sweep trajec-



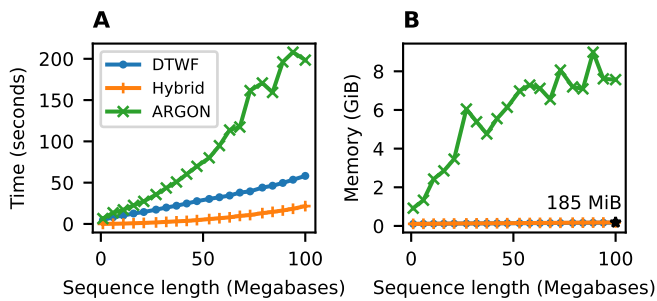
**Figure 6** Comparison of selective sweep simulation performance in `msprime` (`sim_ancestry`) and `discoal` (Intel Xeon Gold 6148 CPU). We report the average CPU time and maximum memory usage when simulating 3 replicates for 100 diploid samples in a model with a single selective sweep in its history, where the beneficial allele had a selection coefficient of  $s = 0.05$ , a per-base recombination rate of  $10^{-8}$ , population size of  $N = 10^4$ , and sequence length varying from 100kb–3000kb.

tories in `msprime` using a jump process approximation to the conditional diffusion of an allele bound for fixation (Coop and Griffiths 2004), as detailed in the appendix. Given a randomly generated allele frequency trajectory, the simulation of a sweep works by assigning lineages to two different structured coalescent “labels”, based on whether they carry the beneficial allele. The allele frequency trajectory determines the relative sizes of the “populations” in these labels over time, and therefore the rates at which various events occur. Common ancestor events can then only merge lineages from *within* a label, but lineages can transfer from one label to the other (i.e., from the advantageous to disadvantageous backgrounds, and vice versa) as a result of recombination events. Once we have reached the end of the simulated trajectory the sweep is complete, and we remove the structured coalescent labels. Simulation may then resume under any other ancestry model.

Fig. 6 compares the performance of `msprime` and `discoal` under a simple sweep model, and shows that `msprime` has far better CPU time and memory performance. Since our implementation uses the abstract label system mentioned above, adding support for similar situations, such as inversions (Peischl *et al.* 2013), should be straightforward.

### Discrete time Wright-Fisher

The coalescent is an idealized model and makes many simplifying assumptions, but it is often surprisingly robust to violations of these assumptions (Wakeley *et al.* 2012). One situation in which the model does break down is the combination of large sample size and long recombining genomes, where the large number of recombination events in the recent past results in more than the biologically possible  $2^t$  ancestors in  $t$  diploid generations (Nelson *et al.* 2020). This pathological behavior results in identity-by-descent, long-range linkage disequilibrium and ancestry patterns deviating from Wright-Fisher expectations, and the bias grows with larger sample sizes (Wakeley *et al.* 2012; Bhaskar *et al.* 2014; Nelson *et al.* 2020). Precisely this problem occurs when simulating modern human datasets, and we have implemented a Discrete Time Wright-Fisher (DTWF) model in `msprime` to address the issue. The DTWF simulates backwards in



**Figure 7** Comparison of Discrete Time Wright-Fisher (DTWF) simulation performance in *msprime* (*sim\_ancestry*) and ARGON (Intel Xeon E5-2680 CPU). We ran simulations with a population size of  $10^4$  and recombination rate of  $10^{-8}$ , with 500 diploid samples and varying sequence length. We report (A) total CPU time and (B) maximum memory usage; each point is the average over 5 replicate simulations. We show observations for ARGON, *msprime*’s DTWF implementation (“DTWF”) and a hybrid simulation of 100 generations of the DTWF followed by the standard coalescent with recombination (“Hybrid”). We ran ARGON with a mutation rate of 0 and with minimum output options, with a goal of measuring only ancestry simulation time. Memory usage for *msprime*’s DTWF and hybrid simulations are very similar.

1 time generation-by-generation so that each gamete has a unique  
 2 diploid parent, and multiple recombinations within a genera-  
 3 tion results in crossover events between the same two parental  
 4 haploid copies. The method is described in detail by Nelson *et al.*  
 5 (2020).

6 Fig. 7 shows that *msprime* simulates the DTWF more quickly  
 7 and requires substantially less memory than ARGON (Palamara  
 8 2016), a specialized DTWF simulator. However, the generation-  
 9 by-generation approach of the DTWF is less efficient than the  
 10 coalescent with recombination when the number of lineages is  
 11 significantly less than the population size (the regime where the  
 12 coalescent is an accurate approximation), which usually hap-  
 13 pens in the quite recent past (Bhaskar *et al.* 2014). We therefore  
 14 support changing the simulation model during a simulation so  
 15 that we can run hybrid simulations, as proposed by Bhaskar *et al.*  
 16 (2014). Any number of different simulation models can be com-  
 17 bined, allowing for the flexible choice of simulation scenarios.  
 18 As the DTWF improves accuracy of genealogical patterns in the  
 19 recent past, we can simulate the recent history using this model  
 20 and then switch to the standard coalescent to more efficiently  
 21 simulate the more ancient history.

### 22 Integration with forward simulators

23 A unique feature of *msprime* is its ability to simulate genetic an-  
 24 cestries by extending an existing partial genetic ancestry. Given  
 25 a tree sequence that is complete up until time  $t$  ago as input  
 26 (where marginal trees may or may not have fully coalesced),  
 27 *msprime* can efficiently obtain the segments of ancestral material  
 28 present at this time, and then run the simulation backwards in  
 29 time from there. This allows a simulated ancestry to be produced  
 30 by any number of different processes across disjoint time slices.  
 31 In practice this feature is used to “complete” forwards-time an-  
 32 cestry simulations (Kelleher *et al.* 2018) that may have not fully  
 33 coalesced. This process (“recapitation”) can be orders of magni-

tude faster than the standard approach of neutral burn-in; see  
 Haller *et al.* (2018) for more details and examples. This interoper-  
 ability between simulators, where a partial ancestry simulation  
 produced by SLiM (Haller and Messer 2019) or fwdpy11 (Thornton  
 2014) can be picked up and completed by another simulator,  
 with complete information retained—at scale—is unprecedented.  
 There may be an opportunity for other forward genetic simula-  
 tors (e.g. Gaynor *et al.* 2021) to leverage the tree sequence data  
 format and associated tools.

### 43 Development model

44 *Msprime* has a large number of features, encompassing the func-  
 45 tionality of several more specialized simulators while maintain-  
 46 ing excellent performance. It is developed by a geographically  
 47 distributed team of volunteers under an open source community  
 48 development model, with a strong emphasis on code quality,  
 49 correctness, good documentation, and inclusive development.  
 50 As in any large code base, unit tests play a key role in ensur-  
 51 ing that new additions behave as expected and *msprime* has an  
 52 extensive suite. These tests are run automatically on different  
 53 operating systems on each pull request (where a contributor pro-  
 54 poses a code change), using standard Continuous Integration  
 55 (CI) methodology. Other CI services check for common errors,  
 56 code formatting issues, and produce reports on the level of test  
 57 coverage for the proposed change.

58 Unit tests are vital for ensuring software quality and correct-  
 59 ness, but they are usually of little value in assessing the statistical  
 60 properties of simulations. To validate the correctness of simu-  
 61 lation output we maintain a suite of statistical tests (as of 1.0.0,  
 62 217 validation tests). These consist of running many replicate  
 63 simulations to check the properties of the output against other  
 64 simulators, and where possible against analytical results. For  
 65 example, simulations of complex demography are validated  
 66 against *ms*, selective sweeps against *discoal*, and Wright-Fisher  
 67 simulations against forwards in time simulations in SLiM. This  
 68 suite of tests is run before every release, to ensure that statistical  
 69 errors have not been introduced.

70 More visibly to the end user, we also have a high standard for  
 71 documentation, with precise, comprehensive, and cross-linked  
 72 documentation that is automatically built from the code base  
 73 and served through the website <https://tskit.dev>. With the  
 74 goal of lowering the entry barrier to new users, we have in-  
 75 vested significant effort in writing examples and introductions,  
 76 and making common tasks discoverable. We also view contri-  
 77 butions to documentation as equally important to the project  
 78 as writing code or designing methods: what use would it be to  
 79 write reliable, stable software if no-one used it?

80 An important goal of *msprime*’s development model is to  
 81 maximize accessibility for prospective users and contributors,  
 82 and to encourage diversity in our community. Gender and  
 83 racial inequality caused by discrimination and marginalization  
 84 is a major problem across the sciences (Wellenreuther and Otto  
 85 2016; Shannon *et al.* 2019) and in open source software develop-  
 86 ment (Trinkenreich *et al.* 2021). Within our field, the contribu-  
 87 tion of women to early computational methods in population gen-  
 88 etics was marginalized (Dung *et al.* 2019), and women continue  
 89 to be under-represented in computational biology (Bonham and  
 90 Stefan 2017). The authorship of our paper reflects these trends,  
 91 with a skew towards men and affiliations in the USA and Europe.  
 92 We know the importance of creating and strengthening networks  
 93 to develop and maintain a diverse community of contributors,  
 94 and we are committed to fostering a supportive and collabora-

1 tive environment that helps to address these inequalities in our  
2 field.

### 3 Discussion

4 The 1.0 release of `msprime` marks a major increase in the breadth  
5 of available features and the potential biological realism of sim-  
6 ulations. These abilities will allow researchers to perform more  
7 robust power analyses, more reliably test new methods, carry  
8 out more reliable inferences, and more thoroughly explore the  
9 properties of theoretical models. Despite this complexity and  
10 generality, `msprime`'s performance is state-of-the-art and all fea-  
11 tures are extensively tested and statistically validated. These  
12 advances have only been possible thanks to a distributed, collab-  
13 orative model of software development, and the work of many  
14 people.

15 Even though simulation has long been a vital tool in popu-  
16 lation genetics, such collaborative software development has  
17 historically been uncommon. A huge proliferation of tools have  
18 been published (the references here are not exhaustive) and only  
19 a small minority of these are actively developed and maintained  
20 today. The ecosystem is highly fragmented, with numerous dif-  
21 ferent ways of specifying parameters and representing results,  
22 and there are significant software quality issues at all stages.  
23 This is unsurprising, since the majority of simulation software  
24 development is performed by students, often without formal  
25 training in software development. The result resembles Hal-  
26 dane's sieve for new mutations: many new pieces of software  
27 stay permanently on a dusty shelf of supplementary materials,  
28 while some of those that prove particularly useful when new  
29 (like dominant alleles) are quickly adopted. Although this has  
30 produced many good tools and enabled decades of research,  
31 it also represents a missed opportunity to invest as a commu-  
32 nity in shared infrastructure and mentorship in good software  
33 development practice.

34 Scientific software is vital apparatus, and must be engineered  
35 to a high quality if we are to trust its results. There is a grow-  
36 ing realization across the sciences (e.g. [Siepel 2019](#); [Harris et al. 2020](#);  
37 [Gardner et al. 2021](#)) that investing in shared community  
38 infrastructure produces better results than a proliferation of in-  
39 dividually maintained tools, allowing scientists to focus on their  
40 specific questions rather than software engineering. `Msprime` 1.0  
41 is the result of such a community process, with features added  
42 by motivated users, taking advantage of the established devel-  
43 opment practices and infrastructure. Software development in  
44 a welcoming community, with mentorship by experienced de-  
45 velopers, is a useful experience for many users. The skills that  
46 contributors learn can lead to greatly increased productivity in  
47 subsequent work (e.g., through more reliable code and better  
48 debugging skills). We hope that users who find that features  
49 they require are missing will continue to contribute to `msprime`,  
50 leading to a community project that is both high quality and  
51 sustainable in the long term.

52 The succinct tree sequence data structure developed for  
53 `msprime` provides a view of not only genetic variation, but also  
54 the genetic ancestry that produced that variation. Recent break-  
55 throughs in methods to infer genetic ancestry in recombining  
56 organisms ([Rasmussen et al. 2014](#); [Kelleher et al. 2019](#); [Speidel et al. 2019](#);  
57 [Wohns et al. 2021](#); [Schaefer et al. 2021](#); [Speidel et al. 2021](#))  
58 have made it possible to estimate such ancestry from real  
59 data at scale for the first time ([Harris 2019](#); [Tang 2019](#)). Given  
60 such inferred ancestry, many exciting applications become possi-  
61 ble. For example, [Osmond and Coop \(2021\)](#) developed a method

to estimate the location of genetic ancestors based on inferred  
trees, and other uses are sure to follow. Since the inferred genetic  
ancestry becomes the input for other downstream inferences, it  
is vitally important that these primary inferences are thoroughly  
validated, with the detailed properties of the inferred ancestries  
cataloged and understood. `Msprime` will continue to be an im-  
portant tool for these inferences and validations, and in this con-  
text the ability to interoperate with other methods—particularly  
forwards simulators—through the succinct tree sequence data  
structure and `tskit` library will be essential.

### Availability

`Msprime` is freely available under the terms of the GNU General  
Public License v3.0, and can be installed from the Python  
Package Index (PyPI) or the `conda-forge` ([conda-forge commu-  
nity 2015](#)) `conda` channel. Development is conducted openly  
on GitHub at <https://github.com/tskit-dev/msprime/>.  
The documentation for `msprime` is available at  
<https://tskit.dev/msprime/docs/>. The source code for  
all the evaluations and figures in this manuscript is available at  
<https://github.com/tskit-dev/msprime-1.0-paper/>.

### Acknowledgments

We acknowledge the contributions of Ivan Krukov who we con-  
sider eligible for authorship, but were unable to contact for  
approval. We would like to thank Iain Mathieson and Alywyn  
Scully for helpful comments on the manuscript.

### Funding

ADK was supported by NIH awards R01GM117241 and  
R01HG010774. AG was supported by NIH award R00HG008696  
to Daniel R Schrider. BE was supported by DFG grant 273887127  
through Priority Programme SPP 1819: Rapid Evolutionary  
Adaptation (grant STE 325/17-2) to Wolfgang Stephan; BE  
would also like to acknowledge funding through The Icelandic  
Research Centre (Rannís) through an Icelandic Research Fund  
Grant of Excellence nr. 185151-051 to Einar Árnason, Katrín  
Halldórsdóttir, Alison Etheridge, Wolfgang Stephan, and BE.  
FB is funded by the Deutsche Forschungsgemeinschaft EXC  
2064/1 – Project number 390727645, and EXC 2124 – Project  
number 390838134. GB and KL are supported by an ERC start-  
ing grant (ModelGenomLand 757648) to KL. Graham Gower  
was supported by a Villum Fonden Young Investigator award  
to Fernando Racimo (project no. 00025300). Gregor Gorjanc  
is supported by the Chancellor's Fellowship of the Univer-  
sity of Edinburgh and the BBSRC grant to The Roslin Institute  
BBS/E/D/30002275. Jere Koskela is supported in part by EP-  
SRC grant EP/R044732/1. Jerome Kelleher is supported by  
the Robertson Foundation. PLR was supported by NIH award  
R01HG010774. SG acknowledges funding from the Canada Re-  
search Chairs Program, from the Canadian Institutes of Health  
Research PJT 173300, and from the Canadian Foundation for  
Innovation.

### Literature Cited

Adrion, J. R., C. B. Cole, N. Dukler, J. G. Galloway, A. L. Glad-  
stein, et al., 2020a A community-maintained standard library  
of population genetic models. *Elife* 9: e54967.  
Adrion, J. R., J. G. Galloway, and A. D. Kern, 2020b Predicting the  
landscape of recombination using deep learning. *Molecular  
biology and evolution* 37: 1790–1808.

Downloaded from <https://academic.oup.com/genetics/advance-article/doi/10.1093/genetics/yab229/6660344> by Edinburgh University user on 06 January 2022

- 1 Arenas, M., 2012 Simulation of molecular data under diverse evolutionary scenarios. *PLoS Computational Biology* **8**: e1002495.
- 2 Arenas, M. and D. Posada, 2007 Recodon: coalescent simulation of coding DNA sequences with recombination, migration and demography. *BMC bioinformatics* **8**: 1–4.
- 3 Árnason, E., 2004 Mitochondrial cytochrome *b* DNA variation in the high-fecundity Atlantic cod: trans-Atlantic clines and shallow gene genealogy. *Genetics* **166**: 1871–1885.
- 4 Barton, N. H., J. Kelleher, and A. M. Etheridge, 2010 A new model for extinction and recolonization in two dimensions: quantifying phylogeography. *Evolution: International journal of organic evolution* **64**: 2701–2715.
- 5 Baumdicker, F. and P. Pfaffelhuber, 2014 The infinitely many genes model with horizontal gene transfer. *Electronic Journal of Probability* **19**: 1–27.
- 6 Beaumont, M. A., W. Zhang, and D. J. Balding, 2002 Approximate Bayesian computation in population genetics. *Genetics* **162**: 2025–2035.
- 7 Becheler, A., C. Coron, and S. Dupas, 2019 The quetzal coalescence template library: A C++ programmers resource for integrating distributional, demographic and coalescent models. *Molecular ecology resources* **19**: 788–793.
- 8 Becheler, A. and L. L. Knowles, 2020 Occupancy spectrum distribution: application for coalescence simulation with generic mergers. *Bioinformatics* **btaa090**.
- 9 Beckenbach, A. T., 1994 Mitochondrial haplotype frequencies in oysters: neutral alternatives to selection models. In *Non-neutral evolution*, pp. 188–198, Springer.
- 10 Bhaskar, A., A. G. Clark, and Y. S. Song, 2014 Distortion of genealogical properties when the sample is very large. *Proceedings of the National Academy of Sciences* **111**: 2385–2390.
- 11 Birkner, M., J. Blath, and B. Eldon, 2013a An ancestral recombination graph for diploid populations with skewed offspring distribution. *Genetics* **193**: 255–290.
- 12 Birkner, M., J. Blath, and B. Eldon, 2013b Statistical properties of the site-frequency spectrum associated with  $\Lambda$ -coalescents. *Genetics* **195**: 1037–1053.
- 13 Birkner, M., J. Blath, M. Möhle, M. Steinrücken, and J. Tams, 2009 A modified lookdown construction for the xi-fleming-viot process with mutation and populations with recurrent bottlenecks. *Alea* **6**: 25–61.
- 14 Birkner, M., H. Liu, and A. Sturm, 2018 Coalescent results for diploid exchangeable population models. *Electronic Journal of Probability* **23**: 1–44.
- 15 Blath, J., M. C. Cronjäger, B. Eldon, and M. Hammer, 2016 The site-frequency spectrum associated with  $\Xi$ -coalescents. *Theoretical Population Biology* **110**: 36–50.
- 16 Blum, M. G. and O. François, 2010 Non-linear regression models for Approximate Bayesian Computation. *Statistics and Computing* **20**: 63–73.
- 17 Bonham, K. S. and M. I. Stefan, 2017 Women are underrepresented in computational biology: An analysis of the scholarly literature in biology, computer science and computational biology. *PLoS computational biology* **13**: e1005134.
- 18 Braverman, J. M., R. R. Hudson, N. L. Kaplan, C. H. Langley, and W. Stephan, 1995 The hitchhiking effect on the site frequency spectrum of DNA polymorphisms. *Genetics* **140**: 783–796.
- 19 Brown, T., X. Didelot, D. J. Wilson, and N. D. Maio, 2016 SimBac: simulation of whole bacterial genomes with homologous recombination. *Microbial Genomics* **2**: 1–6.
- 20 Bunnefeld, L., L. A. F. Frantz, and K. Lohse, 2015 Inferring bottlenecks from genome-wide samples of short sequence blocks. *Genetics* **201**: 1157–1169.
- 21 Bycroft, C., C. Freeman, D. Petkova, G. Band, L. T. Elliott, *et al.*, 2018 The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**: 203–209.
- 22 Cartwright, R. A., 2005 DNA assembly with gaps (Dawg): simulating sequence evolution. *Bioinformatics* **21**: iii31–iii38.
- 23 Carvajal-Rodríguez, A., 2008 Simulation of genomes: a review. *Curr Genomics* **9**: 155.
- 24 Chan, J., V. Perrone, J. P. Spence, P. A. Jenkins, S. Mathieson, *et al.*, 2018 A likelihood-free inference framework for population genetic data using exchangeable neural networks. *Advances in neural information processing systems* **31**: 8594.
- 25 Charlesworth, B. and J. D. Jensen, 2021 Effects of selection at linked sites on patterns of genetic variability. *Annual Review of Ecology, Evolution, and Systematics* **52**: 177–197.
- 26 Charlesworth, B., M. Morgan, and D. Charlesworth, 1993 The effect of deleterious mutations on neutral molecular variation. *Genetics* **134**: 1289–1303.
- 27 Charlesworth, D., B. Charlesworth, and M. Morgan, 1995 The pattern of neutral molecular variation under the background selection model. *Genetics* **141**: 1619–1632.
- 28 Chen, G. K., P. Marjoram, and J. D. Wall, 2009 Fast and flexible simulation of DNA sequence data. *Genome research* **19**: 136–142.
- 29 Chen, H. and K. Chen, 2013 Asymptotic distributions of coalescence times and ancestral lineage numbers for populations with temporally varying size. *Genetics* **194**: 721–736.
- 30 Chen, J.-M., D. N. Cooper, N. Chuzhanova, C. Férec, and G. P. Patrinos, 2007 Gene conversion: mechanisms, evolution and human disease. *Nature Reviews Genetics* **8**: 762–775.
- 31 Chetwynd-Diggles, J. A., B. Eldon, and A. M. Etheridge, 2022 Beta-coalescents when sample size is large. in preparation.
- 32 Chikhi, L., W. Rodríguez, S. Grusea, P. Santos, S. Boitard, *et al.*, 2018 The IICR (inverse instantaneous coalescence rate) as a summary of genomic diversity: insights into demographic inference and model choice. *Heredity* **120**: 13–24.
- 33 conda-forge community, 2015 The conda-forge Project: Community-based Software Distribution Built on the conda Package Format and Ecosystem.
- 34 Coop, G. and R. C. Griffiths, 2004 Ancestral inference on gene trees under selection. *Theoretical population biology* **66**: 219–232.
- 35 Cornuet, J.-M., F. Santos, M. A. Beaumont, C. P. Robert, J.-M. Marin, *et al.*, 2008 Inferring population history with DIY ABC: a user-friendly approach to approximate Bayesian computation. *Bioinformatics* **24**: 2713–2719.
- 36 Csilléry, K., M. G. Blum, O. E. Gaggiotti, and O. François, 2010 Approximate Bayesian computation (ABC) in practice. *Trends in ecology & evolution* **25**: 410–418.
- 37 Csilléry, K., O. François, and M. G. B. Blum, 2012 abc: An R package for approximate Bayesian computation (ABC). *Methods in Ecology and Evolution* **3**: 475–479.
- 38 Dayhoff, M., R. Schwartz, and B. Orcutt, 1978 A model of evolutionary change in proteins. *Atlas of protein sequence and structure* **5**: 345–352.
- 39 De Maio, N., L. Weilguny, C. R. Walker, Y. Turakhia, R. Corbett-Detig, *et al.*, 2021 phastsim: efficient simulation of sequence evolution for pandemic-scale datasets. *bioRxiv*.
- 40 De Maio, N. and D. J. Wilson, 2017 The bacterial sequential markov coalescent. *Genetics* **206**: 333–343.
- 41 Der, R., C. Epstein, and J. B. Plotkin, 2012 Dynamics of neutral and selected alleles when the offspring distribution is skewed.

- 1 Genetics **191**: 1331–1344.
- 2 Desai, M. M., A. M. Walczak, and D. S. Fisher, 2013 Genetic  
3 diversity and the structure of genealogies in rapidly adapting  
4 populations. *Genetics* **193**: 565–585.
- 5 Donnelly, P. and T. G. Kurtz, 1999 Particle representations for  
6 measure-valued population models. *The Annals of Probability*  
7 **27**: 166–205.
- 8 Dung, S. K., A. López, E. L. Barragan, R.-J. Reyes, R. Thu, *et al.*,  
9 2019 Illuminating women’s hidden contribution to historical  
10 theoretical population genetics. *Genetics* **211**: 363–366.
- 11 Durrett, R. and J. Schweinsberg, 2004 Approximating selective  
12 sweeps. *Theoretical population biology* **66**: 129–138.
- 13 Eldon, B. and F. Freund, 2018 Genealogical properties of sub-  
14 samples in highly fecund populations. *Journal of Statistical*  
15 *Physics* **172**: 175–207.
- 16 Eldon, B. and W. Stephan, 2018 Evolution of highly fecund hap-  
17 loid populations. *Theoretical population biology* **119**: 48–56.
- 18 Eldon, B. and J. Wakeley, 2006 Coalescent processes when the  
19 distribution of offspring number among individuals is highly  
20 skewed. *Genetics* **172**: 2621–2633.
- 21 Ethier, S. and R. Griffiths, 1990 On the two-locus sampling dis-  
22 tribution. *Journal of Mathematical Biology* **29**: 131–159.
- 23 Ewing, G. and J. Hermisson, 2010 MSMS: a coalescent simulation  
24 program including recombination, demographic structure,  
25 and selection at a single locus. *Bioinformatics* **26**: 2064–2065.
- 26 Excoffier, L. and M. Foll, 2011 Fastsimcoal: a continuous-time  
27 coalescent simulator of genomic diversity under arbitrarily  
28 complex evolutionary scenarios. *Bioinformatics* **27**: 1332–1334.
- 29 Felsenstein, J. and G. A. Churchill, 1996 A hidden markov model  
30 approach to variation among sites in rate of evolution. *Molec-  
31 ular biology and evolution* **13**: 93–104.
- 32 Flagel, L., Y. Brandvain, and D. R. Schrider, 2019 The unreason-  
33 able effectiveness of convolutional neural networks in popu-  
34 lation genetic inference. *Molecular biology and evolution* **36**:  
35 220–238.
- 36 Fletcher, W. and Z. Yang, 2009 INDELible: a flexible simula-  
37 tor of biological sequence evolution. *Molecular biology and*  
38 *evolution* **26**: 1879–1888.
- 39 Freund, F., 2020 Cannings models, population size changes and  
40 multiple-merger coalescents. *Journal of mathematical biology*  
41 **80**: 1497–1521.
- 42 Galtier, N., F. Depaulis, and N. H. Barton, 2000 Detecting bot-  
43 tlenecks and selective sweeps from DNA sequence polymor-  
44 phism. *Genetics* **155**: 981–987.
- 45 Gardner, P. P., J. M. Paterson, S. R. McGimpsey, F. A. Ghomi,  
46 S. U. Umu, *et al.*, 2021 Sustained software development, not  
47 number of citations or journal choice, is indicative of accurate  
48 bioinformatic software. *bioRxiv* p. 092205.
- 49 Gaynor, R. C., G. Gorjanc, and J. M. Hickey, 2021 AlphaSimR:  
50 An R-package for breeding program simulations. *G3: Genes,*  
51 *Genomes, Genetics* **11**.
- 52 Gillespie, J. H., 2000 Genetic drift in an infinite population: the  
53 pseudohitchhiking model. *Genetics* **155**: 909–919.
- 54 Gladstein, A. L., C. D. Quinto-Cortés, J. L. Pistorius, D. Christy,  
55 L. Gantner, *et al.*, 2018 Simprily: A Python framework to  
56 simplify high-throughput genomic simulations. *SoftwareX*  
57 **7**: 335–340.
- 58 Gower, G., A. P. Ragsdale, *et al.*, 2022 Demes: a standard format  
59 for demographic models. In preparation .
- 60 Griffiths, R. C., 1991 The two-locus ancestral graph. *Lecture*  
61 *Notes-Monograph Series* **18**: 100–117.
- 62 Griffiths, R. C. and P. Marjoram, 1997 An ancestral recombina-  
tion graph. In *Progress in Population Genetics and Human*  
*Evolution, IMA Volumes in Mathematics and its Applications*,  
edited by P. Donnelly and S. Tavaré, volume 87, pp. 257–270,  
Springer-Verlag, Berlin.
- Griffiths, R. C., S. Tavaré, W. F. Bodmer, and P. J. Donnelly, 1994  
Sampling theory for neutral alleles in a varying environment.  
Philosophical Transactions of the Royal Society of London.  
Series B: Biological Sciences **344**: 403–410.
- Guillaume, F. and J. Rougemont, 2006 Nemo: an evolutionary  
and population genetics programming framework. *Bioinforma-  
tics* **22**: 2556–2557.
- Haller, B. C., J. Galloway, J. Kelleher, P. W. Messer, and P. L. Ralph,  
2018 Tree-sequence recording in SLiM opens new horizons  
for forward-time simulation of whole genomes. *Molecular*  
*ecology resources* .
- Haller, B. C. and P. W. Messer, 2019 SLiM 3: forward genetic sim-  
ulations beyond the Wright–Fisher model. *Molecular biology*  
*and evolution* **36**: 632–637.
- Harris, C. R., K. J. Millman, S. J. van der Walt, R. Gommers, P. Vir-  
tanen, *et al.*, 2020 Array programming with numpy. *Nature*  
**585**: 357–362.
- Harris, K., 2019 From a database of genomes to a forest of evolu-  
tionary trees. *Nature genetics* **51**: 1306–1307.
- Hedgcock, D., 1994 Does variance in reproductive success limit  
effective population sizes of marine organisms? *Genetics and*  
*evolution of aquatic organisms* pp. 122–134.
- Hedgcock, D. and A. I. Pudovkin, 2011 Sweepstakes repro-  
ductive success in highly fecund marine fish and shellfish: a  
review and commentary. *Bulletin of Marine Science* **87**: 971–  
1002.
- Hein, J., M. Schierup, and C. Wiuf, 2004 *Gene genealogies, variation*  
*and evolution: a primer in coalescent theory*. Oxford University  
Press, USA.
- Heled, J. and A. J. Drummond, 2009 Bayesian inference of species  
trees from multilocus data. *Molecular biology and evolution*  
**27**: 570–580.
- Hellenthal, G. and M. Stephens, 2007 mshot: modifying Hud-  
son’s ms simulator to incorporate crossover and gene conver-  
sion hotspots. *Bioinformatics* **23**: 520–521.
- Henikoff, S. and J. G. Henikoff, 1992 Amino acid substitution  
matrices from protein blocks. *Proceedings of the National*  
*Academy of Sciences* **89**: 10915–10919.
- Hickerson, M. J., E. Stahl, and N. Takebayashi, 2007 msBayes:  
pipeline for testing comparative phylogeographic histories  
using hierarchical approximate bayesian computation. *BMC*  
*bioinformatics* **8**: 1–7.
- Hoban, S., G. Bertorelle, and O. E. Gaggiotti, 2012 Computer  
simulations: tools for population and evolutionary genetics.  
*Nature Reviews Genetics* **13**: 110–122.
- Hobolth, A. and J. L. Jensen, 2014 Markovian approximation to  
the finite loci coalescent with recombination along multiple  
sequences. *Theoretical population biology* **98**: 48–58.
- Hobolth, A., A. Siri-Jegousse, and M. Bladt, 2019 Phase-type  
distributions in population genetics. *Theoretical population*  
*biology* **127**: 16–32.
- Huang, W., N. Takebayashi, Y. Qi, and M. J. Hickerson, 2011  
MTML-msBayes: approximate Bayesian comparative phylo-  
geographic inference from multiple taxa and multiple loci  
with rate heterogeneity. *BMC bioinformatics* **12**: 1–14.
- Hudson, R. R., 1983a Properties of a neutral allele model with  
intragenic recombination. *Theoretical Population Biology* **23**:  
183–201.

- 1 Hudson, R. R., 1983b Testing the constant-rate neutral allele  
2 model with protein sequence data. *Evolution* **37**: 203–217.
- 3 Hudson, R. R., 1990 Gene genealogies and the coalescent process.  
4 *Oxford Surveys in Evolutionary Biology* **7**: 1–44.
- 5 Hudson, R. R., 2002 Generating samples under a Wright-Fisher  
6 neutral model of genetic variation. *Bioinformatics* **18**: 337–338.
- 7 Irwin, K. K., S. Laurent, S. Matuszewski, S. Vuilleumier, L. Or-  
8 mond, *et al.*, 2016 On the importance of skewed offspring  
9 distributions and background selection in virus population  
10 genetics. *Heredity* **117**: 393–399.
- 11 Johri, P., C. Aquadro, M. Beaumont, B. Charlesworth, L. Excoffier,  
12 *et al.*, 2021 Statistical inference in population genomics .
- 13 Jukes, T. H., C. R. Cantor, *et al.*, 1969 Evolution of protein  
14 molecules. *Mammalian protein metabolism* **3**: 21–132.
- 15 Kamm, J., J. Terhorst, R. Durbin, and Y. S. Song, 2020 Efficiently  
16 inferring the demographic history of many populations with  
17 allele count data. *Journal of the American Statistical Association*  
18 **115**: 1472–1487.
- 19 Kaplan, N. and R. R. Hudson, 1985 The use of sample genealogies  
20 for studying a selectively neutral *m*-loci model with re-  
21 combination. *Theoretical Population Biology* **28**: 382–396.
- 22 Kaplan, N. L., R. R. Hudson, and C. H. Langley, 1989 The “hitch-  
23 hiking effect” revisited. *Genetics* **123**: 887–899.
- 24 Karczewski, K. J., L. C. Francioli, G. Tiao, B. B. Cummings,  
25 J. Alföldi, *et al.*, 2020 The mutational constraint spectrum quan-  
26 tified from variation in 141,456 humans. *Nature* **581**: 434–443.
- 27 Keightley, P. D. and A. Eyre-Walker, 2007 Joint inference of the  
28 distribution of fitness effects of deleterious mutations and  
29 population demography based on nucleotide polymorphism  
30 frequencies. *Genetics* **177**: 2251–2261.
- 31 Kelleher, J., N. H. Barton, and A. M. Etheridge, 2013 Coalescent  
32 simulation in continuous space. *Bioinformatics* **29**: 955–956.
- 33 Kelleher, J., A. M. Etheridge, and N. H. Barton, 2014 Coalescent  
34 simulation in continuous space: Algorithms for large neigh-  
35 bourhood size. *Theoretical population biology* **95**: 13–23.
- 36 Kelleher, J., A. M. Etheridge, and G. McVean, 2016 Efficient coa-  
37 lescent simulation and genealogical analysis for large sample  
38 sizes. *PLoS computational biology* **12**: e1004842.
- 39 Kelleher, J. and K. Lohse, 2020 Coalescent simulation with  
40 msprime. In *Statistical Population Genomics*, edited by J. Y.  
41 Duthel, pp. 191–230, Springer US, New York, NY.
- 42 Kelleher, J., K. R. Thornton, J. Ashander, and P. L. Ralph, 2018  
43 Efficient pedigree recording for fast population genetics simu-  
44 lation. *PLoS Computational Biology* **14**: 1–21.
- 45 Kelleher, J., Y. Wong, A. W. Wohns, C. Fadil, P. K. Albers, *et al.*,  
46 2019 Inferring whole-genome histories in large population  
47 datasets. *Nature Genetics* **51**: 1330–1338.
- 48 Kern, A. D. and D. R. Schrider, 2016 Discoal: flexible coalescent  
49 simulations with selection. *Bioinformatics* **32**: 3839–3841.
- 50 Kim, Y. and W. Stephan, 2002 Detecting a local signature of ge-  
51 netic hitchhiking along a recombining chromosome. *Genetics*  
52 **160**: 765–777.
- 53 Kimura, M., 1980 A simple method for estimating evolutionary  
54 rates of base substitutions through comparative studies of  
55 nucleotide sequences. *Journal of molecular evolution* **16**: 111–  
56 120.
- 57 Kimura, M., 1981 Estimation of evolutionary distances between  
58 homologous nucleotide sequences. *Proceedings of the Na-  
59 tional Academy of Sciences* **78**: 454–458.
- 60 Kingman, J. F., 1982a On the genealogy of large populations.  
61 *Journal of applied probability* **19**: 27–43.
- 62 Kingman, J. F. C., 1982b The coalescent. *Stochastic processes and  
their applications* **13**: 235–248.
- Kluyver, T., B. Ragan-Kelley, F. Pérez, B. Granger, M. Bussonnier,  
*et al.*, 2016 Jupyter notebooks – a publishing format for repro-  
ducible computational workflows. In *Positioning and Power  
in Academic Publishing: Players, Agents and Agendas*, edited by  
F. Loizides and B. Schmidt, pp. 87 – 90, IOS Press.
- Korunes, K. L. and M. A. F. Noor, 2017 Gene conversion and link-  
age: effects on genome evolution and speciation. *Molecular  
Ecology* **26**: 351–364.
- Koskela, J., 2018 Multi-locus data distinguishes between pop-  
ulation growth and multiple merger coalescents. *Statistical  
applications in genetics and molecular biology* **17**.
- Koskela, J. and M. Wilke Berenguer, 2019 Robust model selection  
between population growth and multiple merger coalescents.  
*Mathematical biosciences* **311**: 1–12.
- Kuhner, M. K., J. Yamato, and J. Felsenstein, 2000 Maximum  
likelihood estimation of recombination rates from population  
data. *Genetics* **156**: 1393–1401.
- Lapierre, M., C. Blin, A. Lambert, G. Achaz, and E. P. C. Rocha,  
2016 The impact of selection, gene conversion, and biased sam-  
pling on the assessment of microbial demography. *Molecular  
Biology and Evolution* **33**: 1711–1725.
- Li, H. and R. Durbin, 2011 Inference of human population history  
from individual whole-genome sequences. *Nature* **475**: 493–  
496.
- Li, H. and W. Stephan, 2006 Inferring the demographic history  
and rate of adaptive substitution in *Drosophila*. *PLOS Genet-  
ics* **2**: 1–10.
- Liu, Y., G. Athanasiadis, and M. E. Weale, 2008 A survey of ge-  
netic simulation software for population and epidemiological  
studies. *Human genomics* **3**: 79.
- Lopes, J. S., D. Balding, and M. A. Beaumont, 2009 Popabc: a  
program to infer historical demographic parameters. *Bioinfor-  
matics* **25**: 2747–2749.
- Mailund, T., M. H. Schierup, C. N. Pedersen, P. J. Mechlenborg,  
J. N. Madsen, *et al.*, 2005 CoaSim: a flexible environment  
for simulating genetic data under coalescent models. *BMC  
bioinformatics* **6**: 1–6.
- Marjoram, P. and J. D. Wall, 2006 Fast “coalescent” simulation.  
*BMC Genet* **7**: 16.
- Marth, G. T., E. Czabarka, J. Murvai, and S. T. Sherry, 2004 The  
allele frequency spectrum in genome-wide human variation  
data reveals signals of differential demographic history in  
three large world populations. *Genetics* **166**: 351–372.
- Martin, A. R., C. R. Gignoux, R. K. Walters, G. L. Wojcik, B. M.  
Neale, *et al.*, 2017 Human demographic history impacts ge-  
netic risk prediction across diverse populations. *The American  
Journal of Human Genetics* **100**: 635–649.
- Martin, A. R., C. R. Gignoux, R. K. Walters, G. L. Wojcik,  
B. M. Neale, *et al.*, 2020 Erratum: Human demographic  
history impacts genetic risk prediction across diverse popu-  
lations (the american journal of human genetics (2020)  
107 (4)(583–588),(s000292972030286x),(10.1016/j. ajhg. 2020.08.  
017)). *American journal of human genetics* **107**: 788–789.
- Mathieson, I. and A. Scally, 2020 What is ancestry? *PLoS Genet-  
ics* **16**: e1008624.
- Matuszewski, S., M. E. Hildebrandt, G. Achaz, and J. D. Jensen,  
2018 Coalescent processes with skewed offspring distributions  
and nonequilibrium demography. *Genetics* **208**: 323–338.
- McBroome, J., B. Thornlow, A. S. Hinrichs, N. De Maio, N. Gold-  
man, *et al.*, 2021 A daily-updated database and tools for com-  
prehensive SARS-CoV-2 mutation-annotated trees. *bioRxiv*

- 1  
2 McGill, J. R., E. A. Walkup, and M. K. Kuhner, 2013 GraphML  
3 specializations to codify ancestral recombinant graphs. *Fron*  
4 *Genet* **4**: 146.
- 5 McKenzie, P. F. and D. A. Eaton, 2020 ipcoal: An interactive  
6 Python package for simulating and analyzing genealogies  
7 and sequences on a species tree or network. *Bioinformatics* **36**:  
8 4193–4196.
- 9 McVean, G. A. T. and N. J. Cardin, 2005 Approximating the  
10 coalescent with recombination. *Philos Trans R Soc Lond B Biol*  
11 *Sci* **360**: 1387–1393.
- 12 Miga, K. H., S. Koren, A. Rhie, M. R. Vollger, A. Gershman, *et al.*,  
13 2020 Telomere-to-telomere assembly of a complete human X  
14 chromosome. *Nature* **585**: 79–84.
- 15 Minichiello, M. J. and R. Durbin, 2006 Mapping trait loci by  
16 use of inferred ancestral recombination graphs. *The American*  
17 *Journal of Human Genetics* **79**: 910–922.
- 18 Möhle, M. and S. Sagitov, 2001 A classification of coalescent processes for haploid exchangeable population models. *Annals of Probability* pp. 1547–1562.
- 19 Montinaro, F., V. Pankratov, B. Yelmen, L. Pagani, and M. Mondal, 2020 Revisiting the Out of Africa event with a novel deep learning approach. *bioRxiv*.
- 20 Neher, R. A. and O. Hallatschek, 2013 Genealogies of rapidly adapting populations. *Proceedings of the National Academy of Sciences* **110**: 437–442.
- 21 Nelson, D., J. Kelleher, A. P. Ragsdale, C. Moreau, G. McVean, *et al.*, 2020 Accounting for long-range correlations in genome-wide simulations of large cohorts. *PLoS genetics* **16**: e1008619.
- 22 Nielsen, R., 2000 Estimation of population parameters and recombination rates from single nucleotide polymorphism. *Genetics* **154**: 931–942.
- 23 Osmond, M. and G. Coop, 2021 Estimating dispersal rates and locating genetic ancestors with genome-wide genealogies. *bioRxiv*.
- 24 Palamara, P. F., 2016 ARGON: fast, whole-genome simulation of the discrete time Wright-Fisher process. *Bioinformatics* **32**: 3032–3034.
- 25 Parobek, C. M., F. I. Archer, M. E. DePrenger-Levin, S. M. Hoban, L. Liggins, *et al.*, 2017 skelesim: an extensible, general framework for population genetic simulation in r. *Molecular ecology resources* **17**: 101–109.
- 26 Pavlidis, P., S. Laurent, and W. Stephan, 2010 msABC: a modification of Hudson’s ms to facilitate multi-locus ABC analysis. *Molecular Ecology Resources* **10**: 723–727.
- 27 Peischl, S., E. Koch, R. Guerrero, and M. Kirkpatrick, 2013 A sequential coalescent algorithm for chromosomal inversions. *Heredity* **111**: 200–209.
- 28 Peng, B., H.-S. Chen, L. E. Mechanic, B. Racine, J. Clarke, *et al.*, 2015 Genetic data simulators and their applications: an overview. *Genetic epidemiology* **39**: 2–10.
- 29 Pitman, J., 1999 Coalescents with multiple collisions. *Annals of Probability* pp. 1870–1902.
- 30 Pudlo, P., J. M. Marin, A. Estoup, J. M. Cornuet, M. Gautier, *et al.*, 2016 Reliable ABC model choice via random forests. *Bioinformatics* **32**: 859–866.
- 31 Quinto-Cortés, C. D., A. E. Woerner, J. C. Watkins, and M. F. Hammer, 2018 Modeling SNP array ascertainment with Approximate Bayesian Computation for demographic inference. *Scientific reports* **8**: 1–10.
- 32 Racimo, F., D. Gokhman, M. Fumagalli, A. Ko, T. Hansen, *et al.*, 2017 Archaic adaptive introgression in TBX15/WARS2. *Molecular Biology and Evolution* **34**: 509–524.
- 33 Ragsdale, A. P., D. Nelson, S. Gravel, and J. Kelleher, 2020 Lessons learned from bugs in models of human history. *American Journal of Human Genetics* **107**: 583–588.
- 34 Ralph, P., K. Thornton, and J. Kelleher, 2020 Efficiently summarizing relationships in large samples: a general duality between statistics of genealogies and genomes. *Genetics* **215**: 779–797.
- 35 Rambaut, A. and N. C. Grassly, 1997 Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics* **13**: 235–238.
- 36 Rasmussen, M. D., M. J. Hubisz, I. Gronau, and A. Siepel, 2014 Genome-wide inference of ancestral recombination graphs. *PLoS genetics* **10**: e1004342.
- 37 Raynal, L., J. M. Marin, P. Pudlo, M. Ribatet, C. P. Robert, *et al.*, 2019 ABC random forests for Bayesian parameter inference. *Bioinformatics* **35**: 1720–1728.
- 38 Rivera-Colón, A. G., N. C. Rochette, and J. M. Catchen, 2021 Simulation with RADinitio improves RADseq experimental design and sheds light on sources of missing data. *Molecular ecology resources* **21**: 363–378.
- 39 Rosenzweig, B. K., J. B. Pease, N. J. Besansky, and M. W. Hahn, 2016 Powerful methods for detecting introgressed regions from population genomic data. *Molecular ecology* **25**: 2387–2397.
- 40 Sagitov, S., 1999 The general coalescent with asynchronous mergers of ancestral lines. *Journal of Applied Probability* **36**: 1116–1125.
- 41 Sanchez, T., J. Cury, G. Charpiat, and F. Jay, 2021 Deep learning for population size history inference: Design, comparison and combination with approximate bayesian computation. *Molecular Ecology Resources* **21**: 2645–2660.
- 42 Schaefer, N. K., B. Shapiro, and R. E. Green, 2021 An ancestral recombination graph of human, Neanderthal, and Denisovan genomes. *Science Advances* **7**: eabc0776.
- 43 Schiffels, S. and R. Durbin, 2014 Inferring human population size and separation history from multiple genome sequences. *Nat Genet* **46**: 919–925.
- 44 Schrider, D. R. and A. D. Kern, 2018 Supervised machine learning for population genetics: a new paradigm. *Trends in Genetics* **34**: 301–312.
- 45 Schweinsberg, J., 2000 Coalescents with simultaneous multiple collisions. *Electron Journal of Probability* **5**: 1–50.
- 46 Schweinsberg, J., 2003 Coalescent processes obtained from supercritical Galton–Watson processes. *Stochastic processes and their Applications* **106**: 107–139.
- 47 Schweinsberg, J., 2017 Rigorous results for a population model with selection II: genealogy of the population. *Electronic Journal of Probability* **22**: 1–54.
- 48 Shannon, G., M. Jansen, K. Williams, C. Cáceres, A. Motta, *et al.*, 2019 Gender equality in science, medicine, and global health: where are we at and why does it matter? *The Lancet* **393**: 560–569.
- 49 Sheehan, S., K. Harris, and Y. S. Song, 2013 Estimating variable effective population sizes from multiple genomes: a sequentially markov conditional sampling distribution approach. *Genetics* **194**: 647–662.
- 50 Sheehan, S. and Y. S. Song, 2016 Deep learning for population genetic inference. *PLoS computational biology* **12**: e1004845.
- 51 Shlyakhter, I., P. C. Sabeti, and S. F. Schaffner, 2014 Cosi2: an efficient simulator of exact and approximate coalescent with selection. *Bioinformatics* **30**: 3427–3429.



1 Siepel, A., 2019 Challenges in funding and developing genomic  
2 software: roots and remedies. *Genome Biology* **20**.

3 Speidel, L., L. Cassidy, R. W. Davies, G. Hellenthal, P. Skoglund,  
4 *et al.*, 2021 Inferring population histories for ancient genomes  
5 using genome-wide genealogies. *Molecular Biology and Evo-*  
6 *lution* .

7 Speidel, L., M. Forest, S. Shi, and S. R. Myers, 2019 A method for  
8 genome-wide genealogy estimation for thousands of samples.  
9 *Nature Genetics* **51**: 1321–1329.

10 Spence, J. P. and Y. S. Song, 2019 Inference and analysis of  
11 population-specific fine-scale recombination maps across 26  
12 diverse human populations. *Science Advances* **5**: eaaw9206.

13 Spencer, C. C. and G. Coop, 2004 SelSim: a program to simulate  
14 population genetic data with natural selection and recombina-  
15 tion. *Bioinformatics* **20**: 3673–3675.

16 Spielman, S. J. and C. O. Wilke, 2015 Pyvolve: a flexible Python  
17 module for simulating sequences along phylogenies. *PLoS one*  
18 **10**.

19 Staab, P. R. and D. Metzler, 2016 Coala: an R framework for  
20 coalescent simulation. *Bioinformatics* **32**: 1903–1904.

21 Staab, P. R., S. Zhu, D. Metzler, and G. Lunter, 2015 scrm: Ef-  
22 ficiently simulating long sequences using the approximated  
23 coalescent with recombination. *Bioinformatics* **31**: 1680–1682.

24 Tajima, F., 1983 Evolutionary relationship of DNA sequences in  
25 finite populations. *Genetics* **105**: 437–460.

26 Tang, L., 2019 Genealogy at the genome scale. *Nature methods*  
27 **16**: 1077–1077.

28 Tanjo, T., Y. Kawai, K. Tokunaga, O. Ogasawara, and M. Na-  
29 gasaki, 2021 Practical guide for managing large-scale human  
30 genome data in research. *Journal of Human Genetics* **66**: 39–  
31 52.

32 Tavaré, S. *et al.*, 1986 Some probabilistic and statistical problems  
33 in the analysis of DNA sequences. *Lectures on mathematics*  
34 *in the life sciences* **17**: 57–86.

35 Terasaki Hart, D. E., A. P. Bishop, and I. J. Wang, 2021 Geonomics:  
36 forward-time, spatially explicit, and arbitrarily complex land-  
37 scape genomic simulations. *Molecular Biology and Evolution*  
38 **38**: 4634–4646.

39 Terhorst, J., J. A. Kamm, and Y. S. Song, 2017 Robust and scalable  
40 inference of population history from hundreds of unphased  
41 whole genomes. *Nature genetics* **49**: 303–309.

42 Teshima, K. M. and H. Innan, 2009 mbs: modifying Hudson’s  
43 ms software to generate samples of DNA sequences with a  
44 biallelic site under selection. *BMC Bioinformatics* **10**: 166.

45 Thornton, K. and P. Andolfatto, 2006 Approximate Bayesian  
46 inference reveals evidence for a recent, severe bottleneck in a  
47 Netherlands population of *Drosophila melanogaster*. *Genetics*  
48 **172**: 1607–1619.

49 Thornton, K. R., 2014 A C++ template library for efficient  
50 forward-time population genetic simulation of large popu-  
51 lations. *Genetics* **198**: 157–166.

52 Trinkenreich, B., I. Wiese, A. Sarma, M. Gerosa, and I. Stein-  
53 macher, 2021 Women’s participation in open source software:  
54 A survey of the literature. *arXiv preprint arXiv:2105.08777* .

55 Tskit developers, 2022 Tskit: a portable library for population  
56 scale genealogical analysis. In preparation .

57 Turakhia, Y., B. Thornlow, A. S. Hinrichs, N. De Maio, L. Goza-  
58 shti, *et al.*, 2021 Ultrafast sample placement on existing trees  
59 (USHER) enables real-time phylogenetics for the SARS-CoV-2  
60 pandemic. *Nature Genetics* pp. 1–8.

61 Vendrami, D. L., L. S. Peck, M. S. Clark, B. Eldon, M. Meredith,  
62 *et al.*, 2021 Sweepstake reproductive success and collective

dispersal produce chaotic genetic patchiness in a broadcast  
63 spawner. *Science advances* **7**: eabj4713. 64

Virgoulay, T., F. Rousset, C. Noûs, and R. Leblois, 2021 Gspace:  
65 an exact coalescence simulator of recombining genomes under  
66 isolation by distance. *Bioinformatics* **37**: 3673–3675. 67

Wakeley, J., 2008 *Coalescent theory: an introduction*. Roberts and  
68 Company, Englewood, Colorado. 69

Wakeley, J., L. King, B. S. Low, and S. Ramachandran, 2012 Gene  
70 genealogies within a fixed pedigree, and the robustness of  
71 Kingman’s coalescent. *Genetics* **190**: 1433–1445. 72

Wang, K., I. Mathieson, J. O’Connell, and S. Schiffels, 2020 Track-  
73 ing human population structure through time from whole  
74 genome sequences. *PLoS Genetics* **16**: e1008552. 75

Wang, Y. and B. Rannala, 2008 Bayesian inference of fine-scale  
76 recombination rates using population genomic data. *Philoso-*  
77 *phical Transactions of the Royal Society of London. Series*  
78 *B: Biological Sciences* **363**: 3921–3930. 79

Wang, Y., Y. Zhou, L. Li, X. Chen, Y. Liu, *et al.*, 2014 A new  
80 method for modeling coalescent processes with recombination.  
81 *BMC Bioinformatics* **15**: 273. 82

Wegmann, D., C. Leuenberger, S. Neuenschwander, and L. Ex-  
83 coffier, 2010 ABCtoolbox: a versatile toolkit for approximate  
84 Bayesian computations. *BMC bioinformatics* **11**: 1–7. 85

Wellenreuther, M. and S. Otto, 2016 Women in evolution–  
86 highlighting the changing face of evolutionary biology. *Evolu-*  
87 *tionary Applications* **9**: 3–16. 88

Wilton, P. R., S. Carmi, and A. Hobolth, 2015 The SMC’ is a  
89 highly accurate approximation to the ancestral recombination  
90 graph. *Genetics* **200**: 343–355. 91

Wiuf, C. and J. Hein, 1999a The ancestry of a sample of sequences  
92 subject to recombination. *Genetics* **151**: 1217–1228. 93

Wiuf, C. and J. Hein, 1999b Recombination as a point process  
94 along sequences. *Theoretical Population Biology* **55**: 248–259. 95

Wiuf, C. and J. Hein, 2000 The coalescent with gene conversion.  
96 *Genetics* **155**: 451–462. 97

Wohns, A. W., Y. Wong, B. Jeffery, A. Akbari, S. Mallick, *et al.*,  
98 2021 A unified genealogy of modern and ancient genomes.  
99 *bioRxiv* . 100

Yang, T., H.-W. Deng, and T. Niu, 2014 Critical assessment of  
101 coalescent simulators in modeling recombination hotspots in  
102 genomic sequences. *BMC Bioinformatics* **15**: 3. 103

Yuan, X., D. J. Miller, J. Zhang, D. Herrington, and Y. Wang, 2012  
104 An overview of population genetic data simulation. *Journal*  
105 *of Computational Biology* **19**: 42–54. 106

Zhu, S., J. H. Degnan, S. J. Goldstien, and B. Eldon, 2015 Hybrid-  
107 Lambda: simulation of multiple merger and Kingman gene  
108 genealogies in species networks and species trees. *BMC Bioin-*  
109 *formatics* **16**. 110

## Appendix

### Mutation generation

The algorithm that `msprime` uses to simulate mutations on a tree sequence proceeds in two steps: first, mutations are “placed” on the tree sequence (i.e., sampling their locations in time, along the genome, and on the marginal tree), and then the ancestral and derived alleles of each mutation are generated. All mutation models share the code to place mutations, but choose alleles in different ways.

First, mutations are placed on the tree sequence under an inhomogeneous Poisson model by applying them independently to each edge. If an edge spans a region  $[a, b)$  of the genome

1 and connected parent and child nodes with times  $s < t$ , and  
2 the mutation rate locally is  $\mu$ , then the number of mutations on  
3 the edge is Poisson with mean  $\mu(t - s)(b - a)$ , and each muta-  
4 tion is placed independently at a position chosen uniformly in  
5  $[a, b]$  and a time uniformly in  $[s, t]$ . In a discrete genome, all  
6 positions are integers and so more than one mutation may occur  
7 at the same position on the same edge. Otherwise (i.e., for an  
8 infinite-sites model), positions are rejection sampled to obtain a  
9 unique floating-point number. If an edge spans a region of the  
10 genome with more than one mutation rate, this is done sepa-  
11 rately for each sub-region on which the mutation rate is constant.  
12 Since each edge is processed independently, the algorithm scales  
13 linearly with the number of edges in the tree sequence.

14 Next, alleles are chosen for each mutation. If the site was not  
15 previously mutated, then a new ancestral allele is chosen for the  
16 site, according to an input distribution of ancestral state allele  
17 probabilities. Then, each mutation on the tree is considered  
18 in turn, and a derived allele is randomly chosen based on the  
19 parental allele (which may be the ancestral allele or the derived  
20 allele of a previous mutation). Finally, information about the  
21 mutations are recorded in the site and mutation tables of the tree  
22 sequence.

23 A mutation model must, therefore, provide two things: a way  
24 of choosing an ancestral allele for each new variant site, and a  
25 way of choosing a derived allele given the parental allele at each  
26 mutation. Perhaps the simplest mutation model implemented in  
27 `msprime` is the `InfiniteAlleles` mutation model, which keeps  
28 an internal counter so that the requested alleles are assigned  
29 subsequent (and therefore unique) integers.

30 The distribution of ancestral alleles is used to choose the  
31 allele present at the root of the tree at each mutated site, i.e.,  
32 the `root_distribution`. Mutation models with a finite possible  
33 set of alleles have a natural choice for this distribution—the  
34 *stationary distribution* of the mutation process. (All mutation  
35 models are Markovian, so this may be found as the top left  
36 eigenvector of the mutation matrix.) This is the default in most  
37 models, except, e.g., the `BinaryMutationModel`, whose alleles  
38 are 0 and 1 and always labels the ancestral allele “0”. However,  
39 mutational processes are not in general stationary, so we often  
40 allow a different root distribution to be specified.

41 Since the general algorithm above applies mutations at a  
42 single rate independent of ancestral state, a model in which dif-  
43 ferent alleles mutate at different rates must necessarily produce  
44 some silent mutations, i.e., mutations in which the derived al-  
45 lele is equal to the parental allele. To illustrate this, consider a  
46 mutation model in which *A* or *T* mutates to a randomly chosen  
47 different nucleotide at rate  $\alpha$  and *C* or *G* mutates at rate  $\beta$ , with  
48  $\beta < \alpha$ . To implement this, first place mutations at the largest  
49 total rate, which is  $\alpha$ . Then, at each site, choose an ancestral  
50 allele from the root distribution, and for each mutation, choose  
51 a derived allele as follows: if the parental allele is *A* or *T*, then  
52 choose a random derived allele different to the parental allele; if  
53 the parental allele is *C* or *G*, then choose the derived allele to be  
54 equal to the parent allele with probability  $\beta/(\alpha + \beta)$ , and ran-  
55 domly choose a different nucleotide otherwise. This produces  
56 the correct distribution by Poisson thinning: a Poisson process  
57 with rate  $\alpha$  in which each point is discarded independently with  
58 probability  $\beta/(\alpha + \beta)$  is equivalent to a Poisson process with  
59 rate  $\beta$ . All finite-state models (implemented under the generic  
60 `MatrixMutationModel` class) work in this way: mutations are  
61 placed at the maximum mutation rate, and then some silent  
62 mutations will result.

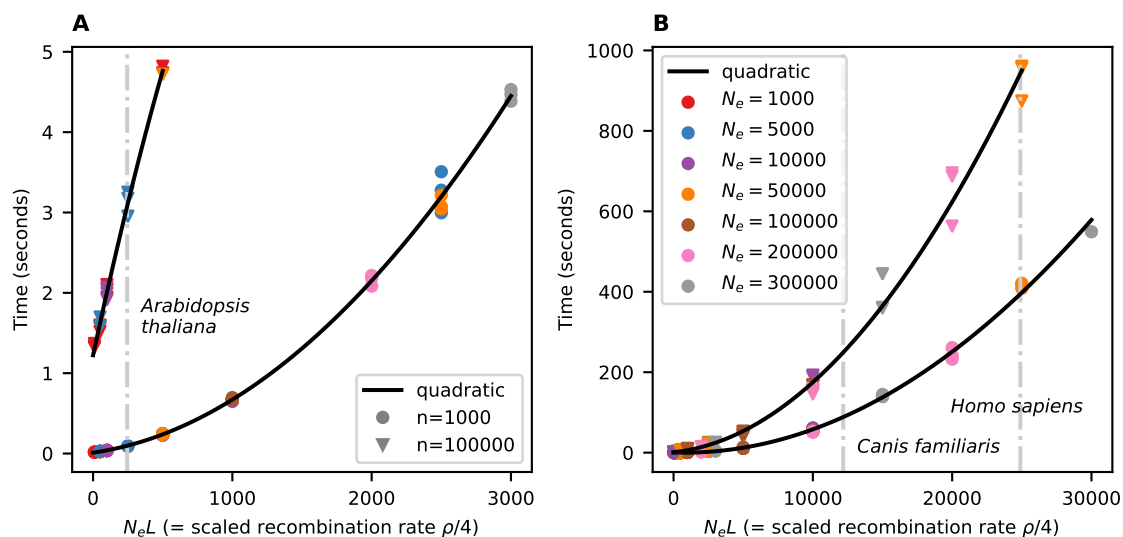
In previous versions of `msprime`, silent mutations were dis-  
allowed, and we could have removed them from the output  
entirely. However, we have chosen to leave them in, so that  
for instance simulating with the HKY mutation model will re-  
sult in silent mutations if not all equilibrium frequencies are  
the same. The presence of silent mutations may at first be sur-  
prising but there is a good reason to leave them in: to allow  
layering of different mutation models. Suppose that we wanted  
to model the mutation process as a mixture of more than one  
model, e.g., Jukes-Cantor mutations at rate  $\mu_1$ , and HKY muta-  
tions occur at rate  $\mu_2$ . Layering multiple calls to `sim_mutations`  
is allowed, so we could first apply mutations with the JC69  
model at rate  $\mu_1$  and then add more with the HKY model at rate  
 $\mu_2$ . However, there is a small statistical problem: suppose that  
after applying Jukes-Cantor mutations we have an  $A \rightarrow C$  mu-  
tation, but then the HKY mutations inserts another mutation  
in the middle, resulting in  $A \rightarrow C \rightarrow C$ . If neither mutation  
model allows silent transitions, then this is clearly not correct,  
i.e., it is not equivalent to a model that simultaneously applies  
the two models. (The impact is small, however, as it only affects  
sites with more than one mutation.) The solution is to make  
the Jukes-Cantor model *state-independent* (also called “parent-  
independent”), by placing mutations at rate  $4/3\mu_1$  and then  
choosing the derived state for each mutation *independently* of  
the parent (so that 1/4 of mutations will be silent). If so—and,  
more generally, if the first mutational process put down is state-  
independent—then the result of sequentially applying the two  
mutation models is equivalent to the simultaneous model. To  
facilitate this, many mutation models have a `state_independent`  
option that increases the number of silent mutations and makes  
the model closer to state-independent.

Silent mutations are fully supported by `tskit`, which cor-  
rectly accounts for their presence when computing statistics and  
performing other operations. For example, silent mutations have  
no effect on calculations of nucleotide site diversity.

### Time complexity of Hudson’s algorithm

As discussed in the section, the time complexity of Hudson’s  
algorithm is predicted to be quadratic in the population scaled re-  
combination rate  $\rho = 4N_eL$  (where  $L$  is the length of the genome  
in units of recombination distance) by Eq. (1). Fig. 8 shows the  
running time for simulations with a variety of population sizes,  
chromosome length and sample sizes, and shows this quadratic  
prediction is well supported by observations (see also Kelleher  
*et al.* 2016, Fig. 2). We also see that the dependence on  $n$  is quite  
weak, since increasing sample size 100-fold only increases run  
time by a factor of 2 or so. However, the  $\log^2 n$  factor implied  
by Eq. (1) (the sum is a harmonic number and can be approxi-  
mated by  $\log n$ ) is not well supported by observed run times (or  
numbers of events) except possibly at very large values of  $\rho$ . It  
therefore appears that a different dependence on  $n$  is required  
to accurately predict simulation time for a given  $\rho$  and  $n$ .

Fig. 8 is a useful yardstick, allowing us to predict how long  
simulations should take for a wide range of species. Taking a  
typical chromosome to be 1 Morgan in length, these plots show,  
roughly, that simulating chromosome-length samples from a  
population of thousands of individuals takes seconds, while  
samples from a population of tens of thousands take minutes.  
Simulating whole chromosomes for many species is very fast,  
with 1000 samples of chromosome 1 for *Arabidopsis thaliana* tak-  
ing less than a second, and a few minutes for dogs and hu-  
mans. However, the dependence on  $\rho$  is quadratic, and if  $\rho$  is



**Figure 8** Running time of `sim_ancestry` for (A) small and (B) larger simulations on an Intel i7-6600U CPU. Each point is the run time of one simulation, for various values of effective population size ( $N_e$ ), chromosome length in Morgans ( $L$ ), and number of diploid samples ( $n$ ). Run time scales quadratically with the product of  $N_e$  and  $L$ , shown on the horizontal axis. For example, (A) shows that 1,000 samples of 1 Morgan-length chromosomes from a population of  $N_e = 2,000$  diploids would take about 2 seconds, and (equivalently) that the same number of 0.01 Morgan segments with  $N_e = 200,000$  would take the same time. Since recombination rate in these simulations was  $10^{-8}$ ,  $L$  is the number of base pairs divided by  $10^8$ . The black lines are quadratic fits separately in each panel and sample size. Vertical gray lines show the approximate values of  $N_e L$  for chromosome 1 in three species, using values from the `stdpopsim` catalog (Adrión *et al.* 2020a).

1 sufficiently large simulations may not be feasible. For exam- 33  
 2 ple, the *Drosophila melanogaster* chromosome 2L is about 23.5Mb 34  
 3 long with an average recombination rate of around  $2.4 \times 10^{-8}$ , 35  
 4 so  $L \approx 0.57$ , and with  $N_e = 1.7 \times 10^6$  (Li and Stephan 2006), 36  
 5  $N_e L \approx 10^6$ , so extrapolating the curve in Fig. 8B predicts that 37  
 6 simulation would require around 177 hours for 1000 samples. 38  
 7 For such large values of  $\rho$  we recommend users consider ap- 39  
 8 proximate simulations. Since `msprime` does not currently have 40  
 9 efficient implementations of approximate coalescent with recom- 41  
 10 bination models, in these cases we recommend using SMC based 42  
 11 methods such as `scrm`, particularly if sample sizes are small. In 43  
 12 practice, to predict the running time of a given simulation in 44  
 13 `msprime`, we recommend that users measure run time in a series 45  
 14 of simulations with short genome lengths and the desired 46  
 15 sample size, and then predict run time by fitting a quadratic 47  
 16 curve to genome length as in Fig. 8. It is important to note that 48  
 17 the quadratic curves in the two panels of Fig. 8 are different, 49  
 18 and predicting the run times of days-long simulations using the 50  
 19 timing of seconds-long runs is unlikely to be very accurate. 51

20 What about simulations with changing population size? To 52  
 21 understand how run time depends on demography it helps 53  
 22 to consider why run time is quadratic in  $\rho$ . At any point in 54  
 23 time, `msprime` must keep track of some number of lineages, each 55  
 24 of which contains some number of chunks of genetic material. 56  
 25 Common ancestor events reduce the number of lineages, and 57  
 26 recombination events increase their number. However, with 58  
 27 long genomes, only a small fraction of the common ancestor 59  
 28 events will involve overlapping segments of ancestry and lead 60  
 29 to coalescence in the marginal trees. Such disjoint segments 61  
 30 are often far apart (on average, about distance  $L/2$ ), and so recom- 62  
 31 bine apart again immediately; it is these large numbers of rapid 63  
 32 and inconsequential events that lead to the quadratic run time. 64

The maximum number of lineages occurs when the increase and 33  
 decrease in numbers of lineages due to common ancestor and 34  
 recombination events balance out. To get an idea of run time we 35  
 can estimate when this balance occurs. Suppose that the maxi- 36  
 mum number of lineages is  $M$ ; at this time the rate of common 37  
 ancestor events is  $M(M-1)/(4N_e)$  and the total rate of recom- 38  
 bination is  $M\ell$ , where  $\ell$  is the mean length of genome carried by 39  
 each lineage (including “trapped” non-ancestral material). At 40  
 the maximum, coalescence and recombination rates are equal, 41  
 so a typical segment of ancestry will spend roughly half its time 42  
 in a lineage with at least one other such segment—and, since 43  
 such lineages carry at least two segments, at most one-third of 44  
 the lineages carry long trapped segments of ancestry. Since the 45  
 maximum number of lineages is reached very quickly (Nelson 46  
*et al.* 2020), this implies that  $\ell \approx L/6$ . Setting the rates of recom- 47  
 bination and common ancestor events to be equal and solving 48  
 for  $M$ , we find that  $M$  is roughly equal to  $LN_e$ . The number of 49  
 lineages then decreases gradually from this maximum on the 50  
 coalescent time scale, and therefore over roughly  $2N_e$  genera- 51  
 tions. Since the total rate of events when the maximum number 52  
 of lineages is present is roughly  $L^2 N_e / 6$ , then the total number 53  
 of events is proportional to  $(LN_e)^2$ —i.e., proportional to  $\rho^2$ . 54

55 What does this tell us about run time for simulating time- 56  
 varying population sizes? Suppose that population size today 57  
 is  $N_1$ , while  $T$  generations ago it was  $N_2$ . Does the run time 58  
 depend more on  $4N_1 L$  or  $4N_2 L$ ? The answer depends on how  $T$  59  
 compares to  $N_1$ : if  $T/N_1 \ll 1$  then the number of extant lineages 60  
 remaining after  $T$  generations is likely to be substantial, and the 61  
 algorithm runtime is primarily determined by  $N_2$ . Conversely, if 62  
 $T/N_1 \gg 1$ , then few extant lineages are likely to remain by time 63  
 $T$  and runtime depends mainly on  $N_1$ . For instance, in many 64  
 agricultural species  $N_1 \approx 100$ , while  $N_2 \approx 10^5$ , and the run time

1 will depend critically on  $T$ —in other words, simulation will be  
 2 quick in a species with a strong domestication bottleneck, and  
 3 slow otherwise.

#### 4 **Selective sweeps model**

Sweep trajectories are generated in `msprime` using a jump process approximation to the conditional diffusion of an allele bound for fixation (Coop and Griffiths 2004). The jump process moves back in time following the beneficial allele frequency,  $p$ , from some initial frequency (e.g.,  $p = 1$ ) back to the origination of the allele at  $p = 1/(2N)$ , tracking time in small increments  $\delta t$ . Then, given the frequency  $p$  at time  $t$ , the frequency  $p'$  at time  $t + \delta t$  is given by

$$p' = \begin{cases} p + \mu(p)\delta t + \sqrt{p(1-p)\delta t} & \text{with probability } 1/2 \\ p + \mu(p)\delta t - \sqrt{p(1-p)\delta t} & \text{with probability } 1/2 \end{cases}$$

5 where

$$\mu(p) = \frac{\alpha p(1-p)}{\tanh(\alpha(1-p))}.$$

6 Here,  $\alpha = 2Ns$  and  $s$  is the fitness advantage in homozygotes.  
 7 This model assumes genic selection (i.e., that the dominance  
 8 coefficient  $h = 0.5$ ), but can be generalized straightforwardly to  
 9 include arbitrary dominance. We can also define trajectories to  
 10 model neutral alleles and soft selective sweeps, which we plan  
 11 as future additions to `msprime`.

#### 12 **Likelihood calculations**

13 We provide two functions to facilitate likelihood-based infer-  
 14 ence. Both are implemented only for the simplest case of the  
 15 standard ARG with a constant population size, and require tree  
 16 sequences compatible with the `record_full_arg` option as their  
 17 arguments.

The `msprime.log_arg_likelihood(ts, r, N)` function re-  
 turns the natural logarithm of the sampling probability of the  
 tree sequence `ts` under the ARG with per-link, per-generation  
 recombination probability `r` and population size `N` (e.g. Kuhner  
 et al. 2000, equation (1)). Specifically, the function returns the  
 logarithm of

$$\left(\frac{1}{2N}\right)^{q_c} \left(\prod_{i:\mathcal{R}} r g_i\right) \exp\left(-\sum_{i=1}^q \left[\frac{1}{2N} \binom{k_i}{2} + r l_i\right] t_i\right),$$

18 where  $t_i$  is the number of generations between the  $(i-1)$ th and  
 19  $i$ th event,  $k_i$  is the number of extant ancestors in that interval,  $l_i$   
 20 is the number of links in that interval that would split ancestral  
 21 material should they recombine,  $q$  is the total number of events  
 22 in the tree sequence `ts`,  $q_c$  is the number of coalescences,  $\mathcal{R}$  is  
 23 the set of indices of time intervals which end in a recombina-  
 24 tion, and  $g_i$  is the corresponding *gap*: the length of contiguous  
 25 non-ancestral material around the link at which the recombina-  
 26 tion in question took place. The gap indicates the number of  
 27 links (or length of genome in a continuous model) at which a  
 28 recombination would result in exactly the observed pattern of  
 29 ancestral material in the ARG. For a continuous model of the  
 30 genome and a recombination in ancestral material, we set  $g_i = 1$   
 31 and interpret the result as a density.

The `msprime.unnormalised_log_mutation_likelihood(ts, m)` function returns the natural logarithm of the probability of the mutations recorded in the tree sequence `ts` given the corresponding ancestry, assuming the infinite sites model,

up to a normalizing constant which depends on the pattern of mutations, but not on the tree sequence or the per-site, per-generation mutation probability  $m$ . Specifically, the function returns the logarithm of

$$e^{-Tm/2} \frac{(Tm/2)^M}{M!} \prod_{\gamma \in \mathcal{M}} \frac{h_\gamma}{T},$$

where  $T$  and  $\mathcal{M}$  are the total branch length and set of mutations in `ts`, respectively, and for a mutation  $\gamma$ ,  $h_\gamma$  is the total branch length on which  $\gamma$  could have arisen while appearing on all of the leaves of `ts` it does, and on no others. Unary nodes on marginal trees arising from the `record_full_arg` option mean that, in general  $h_\gamma$  corresponds to the length of one or more edges.

#### **Multiple merger coalescent model**

Multiple merger coalescents, in which no more than one group of a random number of ancestral lineages may merge into a common ancestor at a given time, are referred to as  $\Lambda$ -coalescents. The rate at which a given group of  $k$  out of a total of  $b$  lineages merges is

$$\lambda_{b,k} = \int_0^1 x^{k-2}(1-x)^{b-k} \Lambda(dx) + a \mathbb{1}_{\{k=2\}}, \quad 2 \leq k \leq b, \quad (2)$$

where  $\mathbb{1}_{\{A\}} := 1$  if  $A$  holds, and zero otherwise,  $a \geq 0$  is a constant, and  $\Lambda$  is a finite measure on the unit interval without an atom at zero (Donnelly and Kurtz 1999; Pitman 1999; Sagitov 1999). There is also a larger class of simultaneous multiple merger coalescents involving simultaneous mergers of distinct groups of lineages into several common ancestors (Schweinsberg 2000). These are commonly referred to as  $\Xi$ -coalescents, and often arise from population models incorporating diploidy or more general polyploidy (Birkner et al. 2013a; Blath et al. 2016). To describe a general  $\Xi$ -coalescent, let  $\Delta$  denote the infinite simplex

$$\Delta := \left\{ (x_1, x_2, \dots) : x_1 \geq x_2 \geq \dots \geq 0, \sum_{j=1}^{\infty} x_j \leq 1 \right\}.$$

The rate of mergers is determined by  $\Xi = \Xi_0 + a\delta_0$ , where  $a \geq 0$  is a constant,  $\delta_0$  is the Dirac delta measure, and  $\Xi_0$  is a finite measure on  $\Delta$  with no atom at  $(0, 0, \dots)$ . For an initial number of blocks  $b \geq 2$  and  $r \in \{1, 2, \dots, b-1\}$ , let  $k_1 \geq 2, \dots, k_r \geq 2$  be the sizes of  $r$  merger events and  $s = b - k_1 - \dots - k_r$  be the number of blocks not participating in any merger. The rate of each possible set of mergers with sizes  $(k_1, \dots, k_r)$  is

$$\begin{aligned} \lambda_{n; k_1, \dots, k_r; s} &= \int_{\Delta} \sum_{\ell=0}^s \sum_{\substack{i_1, \dots, i_{r+\ell}=1 \\ \text{all distinct}}}^{\infty} \binom{s}{\ell} x_{i_1}^{k_1} \dots x_{i_r}^{k_r} x_{i_{r+1}} \dots x_{i_{r+\ell}} \\ &\times \left(1 - \sum_{j=1}^{\infty} x_j\right)^{s-\ell} \frac{1}{\sum_{j=1}^{\infty} x_j^2} \Xi_0(dx) \\ &+ a \mathbb{1}_{\{r=1, k_1=2\}}, \end{aligned}$$

and the number of such  $(k_1, \dots, k_r)$  mergers is

$$\mathcal{N}(b; k_1, \dots, k_r) = \binom{b}{k_1 \dots k_r s} \frac{1}{\prod_{j=2}^b \ell_j!},$$

where  $\ell_j := \#\{i \in \{1, \dots, r\} : k_i = j\}$  is the number of mergers of size  $j \geq 2$  (Schweinsberg 2000).

Viewing coalescent processes strictly as mathematical objects, it is clear that the class of  $\Xi$ -coalescents contains  $\Lambda$ -coalescents as a specific example in which at most one group of lineages can merge at each time, and the class of  $\Lambda$ -coalescents contain the Kingman-coalescent as a special case. However, viewed as limits of ancestral processes derived from specific population models they are not nested. For example, one can obtain  $\Lambda$ -coalescents from haploid population models incorporating sweepstakes reproduction and high fecundity, and  $\Xi$ -coalescents for the same models for diploid populations (Birkner *et al.* 2013a). One should therefore apply the models as appropriate, i.e.  $\Lambda$ -coalescents to haploid (e.g. mtDNA) data, and  $\Xi$ -coalescents to diploid or polyploid (e.g. autosomal) data (Blath *et al.* 2016).

In `msprime` we have incorporated two examples of multiple-merger coalescents. One is a diploid extension (Birkner *et al.* 2013a) of the haploid Moran model adapted to sweepstakes reproduction considered by Eldon and Wakeley (2006). Let  $N$  denote the population size, and take  $\psi \in (0, 1]$  to be fixed. In every generation, with probability  $1 - \varepsilon_N$  a single individual (picked uniformly at random) perishes. With probability  $\varepsilon_N$ ,  $\lfloor \psi N \rfloor$  individuals picked uniformly without replacement perish instead. In either case, a surviving individual picked uniformly at random produces enough offspring to restore the population size back to  $N$ . Taking  $\varepsilon_N = 1/N^\gamma$  for some  $\gamma > 0$ , Eldon and Wakeley (2006) obtain  $\Lambda$ -coalescents for which the  $\Lambda$  measure in (2) is a point mass at  $\psi$ . The simplicity of this model does allow one to obtain some explicit mathematical results (see e.g. Der *et al.* (2012); Eldon and Freund (2018); Freund (2020); Matuszewski *et al.* (2018)), and the model has also been used to simulate gene genealogies within phylogenies (Zhu *et al.* 2015). As well as the haploid model of Eldon and Wakeley (2006), `msprime` provides the diploid version of Birkner *et al.* (2013a), in which individuals perish as above, but replacements are generated by sampling a single pair of diploid individuals as parents, with children sampling one chromosome from each parent. Hence, there are four parent chromosomes involved in each reproduction event, which can lead to up to four simultaneous mergers, giving rise to a  $\Xi$ -coalescent with merger rate

$$\lambda_{b;k_1, \dots, k_r; s}^{\text{Dirac}} = \frac{c\psi^2/4}{1 + c\psi^2/4} \frac{4}{\psi^2} \sum_{\ell=0}^{s \wedge (4-r)} \binom{s}{\ell} (4)_{r+\ell} (1-\psi)^{s-\ell} \times \left(\frac{\psi}{4}\right)^{k_1 + \dots + k_r + \ell} + \frac{\mathbb{1}_{\{r=1, k_1=2\}}}{1 + c\psi^2/4}, \quad (3)$$

The interpretation of (3) is that ‘small’ reproduction events in which two lineages merge occur at rate  $1/(1 + c\psi^2/4)$ , while large reproduction events with the potential to result in simultaneous multiple mergers occur at rate  $(c\psi^2/4)/(1 + c\psi^2/4)$ .

The other multiple merger coalescent model incorporated in `msprime` is the haploid population model considered by Schweinsberg (2003), as well as its diploid extension (Birkner *et al.* 2018). In the haploid version, in each generation of fixed size  $N$ , individuals produce random numbers of juveniles  $(X_1, \dots, X_N)$  independently, each distributed according to a stable law satisfying

$$\lim_{k \rightarrow \infty} Ck^\alpha \mathbb{P}(X \geq k) = 1 \quad (4)$$

with index  $\alpha > 0$ , and where  $C > 0$  is a normalizing constant. If the total number of juveniles  $S_N := X_1 + \dots + X_N$  produced in this way is at least  $N$ , then  $N$  juveniles are sampled uniformly at random without replacement to form the next generation. As

long as  $\mathbb{E}[X_1] > 1$ , one can show that  $\{S_N < N\}$  has exponentially small probability in  $N$ , and does not affect the resulting coalescent as  $N \rightarrow \infty$  (Schweinsberg 2003). If  $\alpha \geq 2$  the ancestral process converges to the Kingman-coalescent; if  $1 \leq \alpha < 2$  the ancestral process converges to a  $\Lambda$ -coalescent with  $\Lambda$  in (2) given by the Beta( $2 - \alpha, \alpha$ ) distribution, i.e.

$$\Lambda(dx) = \mathbb{1}_{\{0 < x \leq 1\}} \frac{1}{B(2 - \alpha, \alpha)} x^{1-\alpha} (1-x)^{\alpha-1} dx, \quad (5)$$

where  $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a+b)$  for  $a, b > 0$  is the beta function (Schweinsberg 2003). This model has been adapted to diploid populations by Birkner *et al.* (2018), and the resulting coalescent is  $\Xi$ -coalescent with merger rate

$$\lambda_{b;k_1, \dots, k_r; s}^{\text{Beta}} = \sum_{\ell=0}^{s \wedge (4-r)} \binom{s}{\ell} \frac{(4)_{r+\ell}}{4^{k+\ell}} \frac{B(k+\ell-\alpha, s-\ell+\alpha)}{B(2-\alpha, \alpha)}, \quad (6)$$

where  $k := k_1 + \dots + k_r$  (Blath *et al.* 2016; Birkner *et al.* 2018). The interpretation of (6) is that the random number of lineages participating in a potential merger is governed by the  $\Lambda$ -coalescent with rate (5), and all participating lineages are randomly allocated into one of four groups corresponding to the four parental chromosomes, giving rise to up to four simultaneous mergers.

The stable law (4) assumes that individuals can produce arbitrarily large numbers of juveniles. Since juveniles are at least fertilized eggs, it may be desirable to suppose that the number of juveniles surviving to reproductive maturity cannot be arbitrarily large. Hence we also consider an adaptation of the Schweinsberg (2003) model, where the random number of juveniles has a deterministic upper bound  $\phi(N)$ , and the distribution of the number of juveniles produced by a given parent (or pair of parents in the diploid case) is

$$\mathbb{P}(X = k) = \mathbb{1}_{\{1 \leq k \leq \phi(N)\}} \frac{\phi(N+1)^\alpha}{\phi(N+1)^\alpha - 1} \left( \frac{1}{k^\alpha} - \frac{1}{(k+1)^\alpha} \right). \quad (7)$$

See Eldon and Stephan (2018) for a related model. One can follow the calculations of Schweinsberg (2003) or Birkner *et al.* (2018) to show that if  $1 < \alpha < 2$  then, recalling that  $k = k_1 + \dots + k_r$ , the merger rate is

$$\lambda_{b;k_1, \dots, k_r; s}^{\text{Beta}, M} = \sum_{\ell=0}^{s \wedge (4-r)} \binom{s}{\ell} \frac{(4)_{r+\ell}}{4^{k+\ell}} \frac{B(M; k+\ell-\alpha, s-\ell+\alpha)}{B(M; 2-\alpha, \alpha)} \quad (8)$$

where  $B(z; a, b) := \int_0^z t^{a-1} (1-t)^{b-1} dt$  for  $a, b > 0$  and  $0 < z \leq 1$  is the incomplete beta function, and

$$M := \lim_{N \rightarrow \infty} \frac{\phi(N)/N}{\phi(N)/N + \mathbb{E}[X_1]} \in (0, 1]$$

(Chetwynd-Diggles *et al.* 2022). In other words, the measure  $\Lambda$  driving the multiple mergers is of the same form as in (5) with  $0 < x \leq M$  in the case  $1 < \alpha < 2$  and  $\lim_{N \rightarrow \infty} \phi(N)/N > 0$ . If  $\alpha \geq 2$  or  $\phi(N)/N \rightarrow 0$  then the ancestral process converges to the Kingman-coalescent (Chetwynd-Diggles *et al.* 2022).