

Demes: a standard format for demographic models

Graham Gower^{1,*}, Aaron P. Ragsdale^{2,*}, Ryan N. Gutenkunst³, Matthew Hartfield⁴, Ekaterina Noskova⁵, Travis J. Struck³, Jerome Kelleher^{6,†,§}, and Kevin R. Thornton^{7,†}

¹Section for Molecular Ecology and Evolution, Globe Institute, University of Copenhagen

²Department of Integrative Biology, University of Wisconsin–Madison

³Department of Molecular and Cellular Biology, University of Arizona

⁴Institute of Evolutionary Biology, The University of Edinburgh

⁵Computer Technologies Laboratory, ITMO University

⁶Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford

⁷Ecology and Evolutionary Biology, University of California at Irvine

*Denotes shared first authorship, listed alphabetically

†Denotes shared senior authorship, listed alphabetically

§Denotes corresponding author

May 31, 2022

Abstract

Understanding the demographic history of populations is a key goal in population genetics, and with improving methods and data, ever more complex models are being proposed and tested. Demographic models of current interest typically consist of a set of discrete populations, their sizes and growth rates, and continuous and pulse migrations between those populations over a number of epochs, which can require dozens of parameters to fully describe. There is currently no standard format to define such models, significantly hampering progress in the field. In particular, the important task of translating the model descriptions in published work into input suitable for population genetic simulators is labor intensive and error prone. We propose the Demes data model and file format, built on widely used technologies, to alleviate these issues. Demes provides a well-defined and unambiguous model of populations and their properties that is straightforward to implement in software, and a text file format that is designed for simplicity and clarity. We provide thoroughly tested implementations of Demes parsers in Python and C, and showcase initial support in several simulators and inference methods.

Introduction

The ever-increasing amount of genetic sequencing data from genetically and geographically diverse species and populations has allowed us to infer complex demography and study life history at fine scales. An integral component to such population genetics studies is simulation. Software to either simulate whole genome sequences (THORNTON, 2014, 2019, STAAB *et al.*, 2015, BAUMDICKER *et al.*, 2022, KELLEHER *et al.*, 2016, HALLER and MESSER, 2019) or informative summary statistics of

diversity (GUTENKUNST *et al.*, 2009, KAMM *et al.*, 2017, JOUGANOUS *et al.*, 2017) have enabled the increasing complexity of genomic studies, with several software packages capable of handling large sample sizes, many interacting populations, and deviations from panmictic random-mating assumptions. This ability to infer and simulate such complex demographic scenarios, however, has highlighted a major shortcoming in community standards: the fragmented landscape of different ways to describe demographic models makes it difficult to compare inferences made by different methods and to reliably simulate from previously inferred models. Inference results are typically reported in publications via a combination of visual depiction, a list of key parameters in tabular form and a discussion within the text. Unfortunately these descriptions are often ambiguous, and implementing the precise model inferred for later simulation is at best tedious and error prone (ADRION *et al.*, 2020, RAGSDALE *et al.*, 2020), and occasionally impossible because of missing information.

Simulation is a core tool in population genetics, and many methods have been developed over the past three decades (CARVAJAL-RODRÍGUEZ, 2008, LIU *et al.*, 2008, ARENAS, 2012, YUAN *et al.*, 2012, HOBAN *et al.*, 2012). Simulations are based on highly idealized population models, and one of the key uses of inferred demographic histories is to make simulations more realistic. Simulation methods take three broad approaches to specifying the demographic model to simulate, using either a command line interface (e.g., HUDSON, 2002, HERNANDEZ, 2008, KERN and SCHRIDER, 2016), a custom input file format (e.g., GUILLAUME and ROUGEMONT, 2006, EXCOFFIER and FOLL, 2011, SHLYAKHTER *et al.*, 2014), or an Application Programming Interface (API) to allow models to be defined programmatically (e.g., THORNTON, 2014, HERNANDEZ and URICCHIO, 2015, KELLEHER *et al.*, 2016, BECHELER *et al.*, 2019, HALLER and MESSER, 2019, THORNTON, 2019, BAUMDICKER *et al.*, 2022). Command line interfaces are a concise way of expressing demographic models, and the syntax defined by `ms` (HUDSON, 2002) is used by several simulators (e.g., EWING and HERMISSON, 2010, CHEN *et al.*, 2009, STAAB *et al.*, 2015). However, this conciseness means that models of even intermediate complexity are difficult for humans to understand, making errors likely. Input parameter file formats for simulators allow the model specification to be less terse and allow for documentation in the form of comments. Several graphical user interfaces and visualization methods have been developed, which greatly facilitate interpretation (MAILUND *et al.*, 2005, ANTAO *et al.*, 2007, PARREIRA *et al.*, 2009, EWING and HERMISSON, 2010, PAROBEK *et al.*, 2017, ZHOU *et al.*, 2018). However, these methods currently have little traction as they are all either directly coupled to an internal simulation method or to the syntax of a specific simulator. There is currently no way in which demographic models inferred by different packages can be simulated or visualized by downstream software.

Here we present “Demes”, a data model and file format specification for complex demographic models developed by the PopSim Consortium (ADRION *et al.*, 2020). The Demes data model precisely defines the sizes and relationships of populations, and it provides a way to explicitly encode the information relevant to demography while avoiding repetition. This data model is implemented in the widely used YAML format (BEN-KIKI *et al.*, 2009), which provides a good balance between human and machine readability. The specification precisely defines the required behavior of implementations, ensuring that there is no ambiguity of interpretation, and includes both a reference implementation and an extensive suite of test examples and their expected output. The initial software ecosystem includes high-quality Python and C parser implementations, as well as utilities for verification and visualization of Demes models, and has been implemented in several popular inference and simulation methods (Table 1). We hope that this data model and file format will be widely adopted by the community, such that users can expect to simulate directly from

inferred models with little to no programming effort.

Demes

The design of Demes is a balance between two partially competing requirements: that (a) models should be easy for humans to understand and manipulate; and (b) software processing Demes models should be provided with an unambiguous representation that is straightforward to process. For efficiency of understanding and avoidance of model specification error, we require a data representation without redundancy (i.e., repetition of values). However, for the simplicity of software working with the Demes model (and the avoidance of programming error, or divergence in interpretations of the specification) it is preferable to have an explicit representation, in which all relevant values are readily available. Thus, Demes is composed of three entities: the Human Data Model (HDM) designed for human readability; the Machine Data Model (MDM) designed for programmatic input and processing; and the parser, which is responsible for transforming the former into the latter.

Below we provide a brief overview of the population genetics models that Demes supports and the components of the Demes infrastructure. We then highlight an example usage, including visualization and simulation of a multi-population demographic model. More detailed descriptions of the population genetics model, parsers, scope, and additional examples are given in the Appendix. Complete technical details of the MDM and HDM, and the responsibilities of the parser are provided in the online Demes specification (<https://popsim-consortium.github.io/demes-spec-docs/>).

Population genetics model

Demographic models consist of one or more populations (or “demes”) defined by their size histories and the time intervals of their existence. Individuals can move between populations based on their ancestor-descendant relationships or by continuous or discrete migration events. The demographic model itself omits information about genome biology, such as genome size and architecture, ploidy, and mutation and recombination rates. It also omits any information about selection operating within populations. Within populations, allele-frequency dynamics are assumed to follow the Wright-Fisher model, so that generations are discrete and non-overlapping and population sizes are independent of mean fitness (i.e., “soft” selection (CHRISTIANSEN, 1975)). While these restrictions preclude some population-genetic scenarios of interest, such as continuous spatial structure (WRIGHT, 1943, BARTON *et al.*, 2002, 2010, RINGBAUER *et al.*, 2017, BATTEY *et al.*, 2020) or stochastic population size dynamics (NUNNEY and CAMPBELL, 1993, ORR and UNCKLESS, 2014), many current forwards- and backwards-in-time simulators make these same assumptions. The basic assumptions of discrete populations connected by instantaneous or continuous migrations are also shared by many inference methods (e.g., GUTENKUNST *et al.*, 2009, GRAVEL, 2012, KAMM *et al.*, 2017, JOUGANOUS *et al.*, 2017, RAGSDALE and GRAVEL, 2019, EXCOFFIER *et al.*, 2021). Demes therefore serves as “middleware” between inference methods and simulation software, capturing these common assumptions. Appendix A1 provides more details on the population genetics models supported by Demes.

Human Data Model

The Demes Human Data Model (HDM) is focused on efficient human understanding and avoiding errors. We have adopted the widely used YAML format (BEN-KIKI *et al.*, 2009) as the primary

Software infrastructure	
<code>demes-python</code>	A Python library for loading, saving, and working with Demes models. Includes support for converting to and from <code>ms</code> (HUDSON, 2002) (https://github.com/popsim-consortium/demes-python).
<code>demes-c</code>	A C library for parsing Demes YAML descriptions. (https://github.com/grahamgower/demes-c).
<code>demesdraw</code>	A Python library for visualizing Demes models (https://github.com/grahamgower/demesdraw).

Methods using Demes as input/output format	
<code>dadi</code>	Optimizes parameters in models of demographic history and distributions of fitness effects using SFS (GUTENKUNST <i>et al.</i> , 2009). Can simulate SFS from Demes models.
<code>demes-slim</code>	Loads Demes models into the SLiM forward simulator (HALLER and MESSER, 2019).
<code>fdpy11</code>	Simulates the Wright-Fisher model forward in time (THORNTON, 2014, 2019). Demes is the preferred format for specifying a demographic model.
<code>GADMA</code>	Infers models of demographic history (NOSKOVA <i>et al.</i> , 2020). Outputs Demes models and visualizations.
<code>moments</code>	Optimizes parameters in models of demographic history using SFS and linkage disequilibrium statistics (JOUGANOUS <i>et al.</i> , 2017, RAGSDALE and GRAVEL, 2019). Models to be optimized can be specified in Demes.
<code>msprime</code>	Simulates population genetic models using tree sequences (KELLEHER <i>et al.</i> , 2016, KELLEHER and LOHSE, 2020, BAUMDICKER <i>et al.</i> , 2022). Demographic history models can be specified using Demes.

Table 1: **Software support for Demes.** We have included software infrastructure developed for working with Demes models (such as parsing, validation, and visualization) as well as downstream software that implement the specification, at the time of writing.

interface for writing and interchanging demographic models (see Appendix A2 for details about our choice of YAML). Demographic models written in YAML format provide information about global features of the model (such as time units and generation times), populations (as “demes”) and their existence intervals (as “epochs”), and gene flow between populations (as continuous “migrations” or instantaneous “pulse” events) (e.g., Figures 1 and A1). The HDM encourages human understanding by avoiding redundancy in the description where possible and by providing a mechanism for specifying default values that are inherited hierarchically.

Parsers

While the HDM is designed for human readability and conciseness, the underlying data model suitable for software implementation (the Machine Data Model, or MDM) is redundant and exhaustive. Translation from the HDM to the MDM requires resolving hierarchically-defined default values and verifying relationships between populations and the validity of specified parameter values. Because this translation and validation requires significant programming effort, we define a standard software entity as part of the specification to perform this task (the parser), which is intended to be shared by programs that support Demes as input. By maintaining high-quality Demes parsers available as libraries, we ensure consistency across simulation and inference software. In addition to a reference implementation written in Python, we provide high-quality parser implementations in the Python and C languages (Table 1, as discussed in Appendix A3).

The scope of Demes

The Demes specification is static by design—we wish to unambiguously describe a demographic model with a concrete set of parameters. This simplicity means that we cannot directly specify parameter distributions or estimated confidence intervals for those parameters. While it is not difficult to imagine extending the specification in ways that would allow this (Appendix A4), it is not immediately clear that the benefits are worth the increased parser complexity. It is important to note, however, that Demes may be used in applications that include additional population genetic processes outside of what is explicitly modeled in the specification, such as interpreting population sizes as carrying capacities, implementations of hard selection, or layering more complicated mating or spatial structure. The Demes specification therefore provides a basic model that can be elaborated on where necessary.

Application: simulation using Demes

Here, we highlight the interaction between Demes and other software, including simulation and model illustration tools. Demes allows us to specify a demographic model which can be used as the input for a growing number of simulation packages (Table 1). We implemented the human two-population demographic model from TENNESSEN *et al.* (2012) inferred from European and African-American sequencing data. This model (shown in Demes format in Figure A2) is parameterized by an ancestral population with an ancient growth, divergence into “AFR” and “EUR” that each have multiple-epoch size histories, and multiple epochs of continuous migration between the two branches (illustrated using `demesdraw` in Figure 1A). The large final sizes ($\approx 500,000$ individuals each) are one to three orders of magnitude larger than ancestral population sizes, reflecting the recent explosive population size increase in humans.

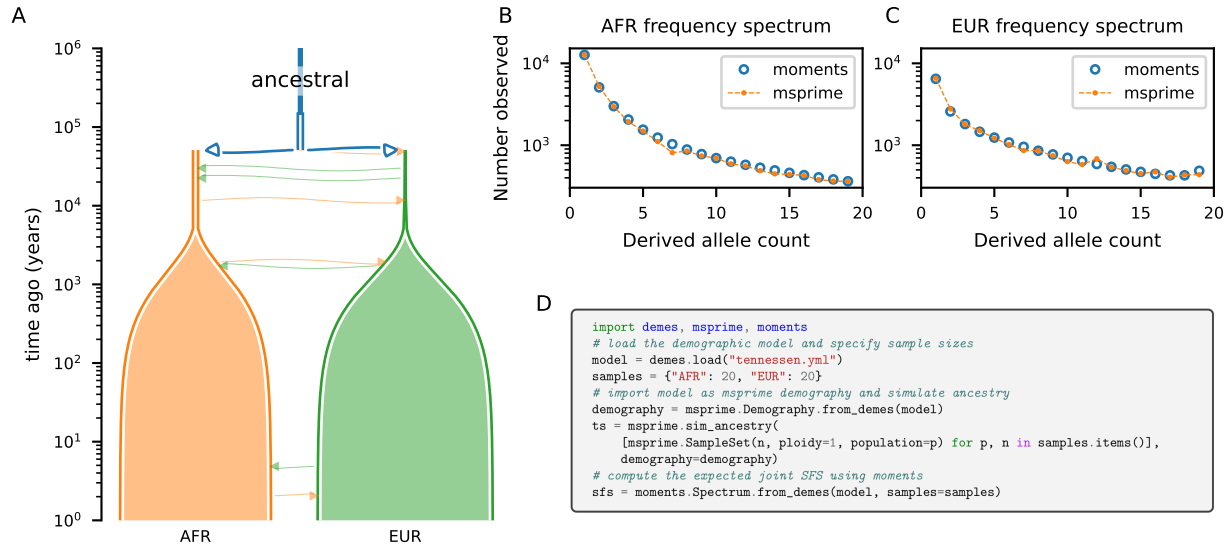


Figure 1: Illustration and simulation using Demes. (A) Using an inferred demographic model from TENNESSEN *et al.* (2012) specified as a YAML file in Demes format (Figure A2), we used `demesdraw` to visualize the demographic model (note the recent exponential growth resulting in present-day population sizes that greatly exceed those in the past). We then used `msprime` to simulate genomic data for 20 genome copies sampled from the two contemporary populations, and we used `moments` to compute the expected joint site-frequency spectrum for the same sample sizes (Figure A3). (B, C) We compared the single-population SFS in each population, showing agreement between the simulation methods. (D) Python code snippets of the interactions between `demes` and the simulation software. An extended script to compute the SFS shown in (B) and (C) is given in Figure A3.

We used this model to simulate 20 haploid genome copies from EUR and AFR at time zero (i.e., present day) to obtain the joint site-frequency spectrum (SFS), a summary of observed allele frequencies widely used in evolutionary inference (BUSTAMANTE *et al.*, 2001, GUTENKUNST *et al.*, 2009, TENNESSEN *et al.*, 2012, JOUGANOUS *et al.*, 2017, KAMM *et al.*, 2017, KIM *et al.*, 2017). The Demes model (Figures 1A and A2) was provided as the input demography to `msprime` (BAUMDICKER *et al.*, 2022) to simulate a large recombining region under the mutation rate assumed in TENNESSEN *et al.* (2012), and we computed the observed SFS using `tskit` (RALPH *et al.*, 2020). Using the same Demes model as input to `moments` (JOUGANOUS *et al.*, 2017), we computed the expectation of the joint SFS and compared to the `msprime` simulated data (Figure 1B,C). Figure 1D shows the code required to run the simulations in `msprime` and `moments`, and demonstrates that precisely the same input model, without modification, was provided to both packages. Such interoperability is a major gain for researchers, which we hope will become the expected norm as more packages adopt the Demes format.

Discussion

Stable and healthy software ecosystems require standard interchange formats, allowing for the development of high-quality and long-lasting tools that produce and consume the standard. Demographic models are a key part of population genetics research, and to date the transfer of inferred models to downstream simulations has been *ad-hoc*, and conversions between the many different ways of expressing such models is both labor intensive and fraught with errors. The proposed Demes standard is an attempt to bridge this gap between inference and simulation, and also to provide the foundations for a sustainable ecosystem of tools built around this data model. Table 1 shows some initial infrastructure that we have built as part of developing Demes, but many other useful tools can be envisaged that consume, transform, or produce this format.

Reproducibility is a significant problem throughout the sciences (BAKER, 2016), and various measures have been proposed to increase the likelihood of researchers being able to replicate results in the literature (MUNAFÒ *et al.*, 2017). The most basic requirement for reproducibility is that we must be able to state precisely what the result in question is. The lack of standardization in how complex demographic models are communicated today, and the lack of precision in the published model descriptions means that it is difficult to replicate analyses, or reproduce those models for later simulation. Thus, we hope that the Demes standard introduced here will be widely adopted by simulation and inference methods and be used for reporting results in publications, either as supplemental material or uploaded to a data repository.

References

- ADRION, J. R., C. B. COLE, N. DUKLER, J. G. GALLOWAY, A. L. GLADSTEIN, *et al.*, 2020 A community-maintained standard library of population genetic models. *eLife* **9**: e54967.
- ANTAO, T., A. BEJA-PEREIRA, and G. LUIKART, 2007 MODELER4SIMCOAL2: A user-friendly, extensible modeler of demography and linked loci for coalescent simulations. *Bioinformatics* **23**: 1848–1850.
- ARENAS, M., 2012 Simulation of molecular data under diverse evolutionary scenarios. *PLoS Computational Biology* **8**: e1002495.
- BAKER, M., 2016 1,500 scientists lift the lid on reproducibility. *Nature News* **533**: 452.
- BARTON, N. H., F. DEPAULIS, and A. M. ETHERIDGE, 2002 Neutral evolution in spatially continuous populations. *Theoretical Population Biology* **61**: 31–48.
- 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.
- BATTEY, C., P. L. RALPH, and A. D. KERN, 2020 Space is the place: effects of continuous spatial structure on analysis of population genetic data. *Genetics* **215**: 193–214.
- BAUMDICKER, F., G. BISSCHOP, D. GOLDSTEIN, G. GOWER, A. P. RAGSDALE, *et al.*, 2022 Efficient ancestry and mutation simulation with msprime 1.0. *Genetics* **220**: iyab229.

- BEAUMONT, M. A., W. ZHANG, and D. J. BALDING, 2002 Approximate Bayesian computation in population genetics. *Genetics* **162**: 2025–2026.
- 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.
- BEN-KIKI, O., C. EVANS, and B. INGERSON, 2009 YAML ain't markup language (yaml™) version 1.1. Working Draft 2008-05 **11**.
- BRAY, T., 2017 The JavaScript Object Notation (JSON) Data Interchange Format. RFC 8259.
- BÜRGER, R., 2000 *The Mathematical Theory of Selection, Recombination, and Mutation*. Wiley.
- BUSTAMANTE, C. D., J. WAKELEY, S. SAWYER, and D. L. HARTL, 2001 Directional selection and the site-frequency spectrum. *Genetics* **159**: 1779–1788.
- CARVAJAL-RODRÍGUEZ, A., 2008 Simulation of genomes: a review. *Current Genomics* **9**: 155.
- CHEN, G. K., P. MARJORAM, and J. D. WALL, 2009 Fast and flexible simulation of DNA sequence data. *Genome Research* **19**: 136–142.
- CHRISTIANSEN, F. B., 1975 Hard and soft selection in a subdivided population. *The American Naturalist* **109**: 11–16.
- CROW, J. F., and M. KIMURA, 1970 *An introduction to mathematical population genetics theory*. Alpha Editions.
- EWING, G., and J. HERMISSON, 2010 MSMS: a coalescent simulation program including recombination, demographic structure, and selection at a single locus. *Bioinformatics* **26**: 2064–2065.
- EXCOFFIER, L., and M. FOLL, 2011 Fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics* **27**: 1332–1334.
- EXCOFFIER, L., N. MARCHI, D. A. MARQUES, R. MATTHEY-DORET, A. GOUY, *et al.*, 2021 fastsimcoal2: demographic inference under complex evolutionary scenarios. *Bioinformatics* **37**: 4882–4885.
- GILMOUR, J. S. L., and J. W. GREGOR, 1939 Demes: A suggested new terminology. *Nature* **144**: 333–333.
- GILMOUR, J. S. L., and J. HESLOP-HARRISON, 1955 The deme terminology and the units of micro-evolutionary change. *Genetica* **27**: 147–161.
- GRAVEL, S., 2012 Population genetics models of local ancestry. *Genetics* **191**: 607–619.
- GUILLAUME, F., and J. ROUGEMONT, 2006 Nemo: an evolutionary and population genetics programming framework. *Bioinformatics* **22**: 2556–2557.
- GUTENKUNST, R. N., R. D. HERNANDEZ, S. H. WILLIAMSON, and C. D. BUSTAMANTE, 2009 Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genetics* **5**: e1000695.

- HALLER, B. C., and P. W. MESSER, 2019 SLiM 3: forward genetic simulations beyond the Wright–Fisher model. *Molecular Biology and Evolution* **36**: 632–637.
- HARTFIELD, M., S. I. WRIGHT, and A. F. AGRAWAL, 2016 Coalescent Times and Patterns of Genetic Diversity in Species with Facultative Sex: Effects of Gene Conversion, Population Structure, and Heterogeneity. *Genetics* **202**: 297–312.
- HERNANDEZ, R. D., 2008 A flexible forward simulator for populations subject to selection and demography. *Bioinformatics* **24**: 2786–2787.
- HERNANDEZ, R. D., and L. H. URICCHIO, 2015 SFS_code: More Efficient and Flexible Forward Simulations. Technical report, bioRxiv.
- HOBAN, S., G. BERTORELLE, and O. E. GAGGIOTTI, 2012 Computer simulations: tools for population and evolutionary genetics. *Nature Reviews Genetics* **13**: 110–122.
- HUDSON, R. R., 1983 Testing the constant-rate neutral allele model with protein sequence data. *Evolution* : 203–217.
- HUDSON, R. R., 2002 Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics* **18**: 337–338.
- JOUGANOUS, J., W. LONG, A. P. RAGSDALE, and S. GRAVEL, 2017 Inferring the joint demographic history of multiple populations: beyond the diffusion approximation. *Genetics* **206**: 1549–1567.
- KAMM, J. A., J. TERHORST, and Y. S. SONG, 2017 Efficient computation of the joint sample frequency spectra for multiple populations. *Journal of Computational and Graphical Statistics* **26**: 182–194.
- KELLEHER, J., A. M. ETHERIDGE, and G. MCVEAN, 2016 Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Computational Biology* **12**: e1004842.
- KELLEHER, J., and K. LOHSE, 2020 Coalescent simulation with msprime. In J. Y. Dutheil, editor, *Statistical Population Genomics*. Springer US, New York, NY, 191–230.
- KERN, A. D., and D. R. SCHRIDER, 2016 Discoal: flexible coalescent simulations with selection. *Bioinformatics* **32**: 3839–3841.
- KIM, B. Y., C. D. HUBER, and K. E. LOHMUELLER, 2017 Inference of the distribution of selection coefficients for new nonsynonymous mutations using large samples. *Genetics* **206**: 345–361.
- LIU, Y., G. ATHANASIADIS, and M. E. WEALE, 2008 A survey of genetic simulation software for population and epidemiological studies. *Human Genomics* **3**: 79.
- 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.
- MUNAFÒ, M. R., B. A. NOSEK, D. V. BISHOP, K. S. BUTTON, C. D. CHAMBERS, *et al.*, 2017 A manifesto for reproducible science. *Nature Human Behaviour* **1**: 1–9.

- NORDBORG, M., and P. DONNELLY, 1997 The coalescent process with selfing. *Genetics* **146**: 1185–1195.
- NOSKOVA, E., V. ULYANTSEV, K.-P. KOEPFLI, S. J. O'BRIEN, and P. DOBRYNIN, 2020 GADMA: Genetic algorithm for inferring demographic history of multiple populations from allele frequency spectrum data. *GigaScience* **9**: giaa005.
- NUNNEY, L., and K. A. CAMPBELL, 1993 Assessing minimum viable population size: demography meets population genetics. *Trends in Ecology & Evolution* **8**: 234–239.
- ORR, H. A., and R. L. UNCKLESS, 2014 The population genetics of evolutionary rescue. *PLoS Genetics* **10**: e1004551.
- 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.
- PARREIRA, B., M. TRUSSART, V. SOUSA, R. HUDSON, and L. CHIKHI, 2009 SPAMs: A user-friendly software to simulate population genetics data under complex demographic models. *Molecular Ecology Resources* **9**: 749–753.
- RAGSDALE, A. P., and S. GRAVEL, 2019 Models of archaic admixture and recent history from two-locus statistics. *PLoS Genetics* **15**: e1008204.
- 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.
- 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.
- RINGBAUER, H., G. COOP, and N. H. BARTON, 2017 Inferring recent demography from isolation by distance of long shared sequence blocks. *Genetics* **205**: 1335–1351.
- SHLYAKHTER, I., P. C. SABETI, and S. F. SCHAFFNER, 2014 C_{osi}2: an efficient simulator of exact and approximate coalescent with selection. *Bioinformatics* **30**: 3427–3429.
- STAAB, P. R., S. ZHU, D. METZLER, and G. LUNTER, 2015 scrm: Efficiently simulating long sequences using the approximated coalescent with recombination. *Bioinformatics* **31**: 1680–1682.
- TAJIMA, F., 1983 Evolutionary relationship of DNA sequences in finite populations. *Genetics* **105**: 437–460.
- TENNESSEN, J. A., A. W. BIGHAM, T. D. O'CONNOR, W. FU, E. E. KENNY, *et al.*, 2012 Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science* **337**: 64–69.
- THORNTON, K. R., 2014 A C++ template library for efficient forward-time population genetic simulation of large populations. *Genetics* **198**: 157–166.
- THORNTON, K. R., 2019 Polygenic adaptation to an environmental shift: Temporal dynamics of variation under gaussian stabilizing selection and additive effects on a single trait. *Genetics* **213**: 1513–1530.

- WAKELEY, J., 2008 *Coalescent Theory: An Introduction*. W. H. Freeman.
- WRIGHT, A., H. ANDREWS, B. HUTTON, and G. DENNIS, 2020 JSON schema: A media type for describing JSON documents.
- WRIGHT, S., 1943 Isolation by distance. *Genetics* **28**: 114.
- YUAN, X., D. J. MILLER, J. ZHANG, D. HERRINGTON, and Y. WANG, 2012 An overview of population genetic data simulation. *J Comput Biol* **19**: 42–54.
- ZHOU, Y., X. TIAN, B. L. BROWNING, and S. R. BROWNING, 2018 POPdemog: visualizing population demographic history from simulation scripts. *Bioinformatics* **34**: 2854–2855.

Appendix

The Demes specification is a formal data model for describing the properties of populations over time, along with some metadata and provenance information. The data model is based on the ubiquitous JSON (BRAY, 2017) standard, and formally defined using JSON Schema (WRIGHT *et al.*, 2020). Along with the schema, full technical details of the of the model are provided in the on-line specification document (<https://popsim-consortium.github.io/demes-spec-docs/>). This specification rigorously defines the data model, fully describing the entities and their relationships, and the required behavior of implementations. Since the online specification is definitive, we will not recapitulate the details here, but instead focus on the high level properties of the model and the rationale behind key design decisions.

Below we detail the population genetics model in Appendix A1, the high-level, human-readable model specification in Appendix A2, demographic model parsers and their intended behavior in Appendix A3, and the intended scope of Demes in Appendix A4. The extended figures in Appendix A5 provide additional examples and illustrations of Demes usage.

A1 Population genetics model details

In Demes, demographic models consist of one or more interacting populations, or “demes”, understood to be a collection of individuals that can be conveniently modeled using a defined set of rules and parameters (GILMOUR and GREGOR, 1939, GILMOUR and HESLOP-HARRISON, 1955). To avoid confusion with the name of the specification itself we will use the term “population” in this discussion, with the understanding that the terms are interchangeable. A population is defined as some collection of individuals that exists for some period of time, and has a well-defined size (i.e., number of individuals) during that time period. Individuals can move between populations either according to their ancestor-descendant relationships or through processes involving migrations. Few other properties of the populations are specified in the model: we are concerned primarily with defining the populations, their sizes, and the movement of individuals between those populations.

A1.1 Time units

Population and event times are written as units in the past, so that time zero corresponds to the final generation or “now”, and event times in the past are values greater than zero with larger values

corresponding to times in the more distant past. By having time units increase into the past, we avoid the need to choose an arbitrary point in history as “time zero”. A natural specification for time units is in generations, although other time units are permitted, such as years, accompanied by the generation time.

A1.2 Sizes and epochs

Population sizes are given as numbers of individuals, and details such as ploidy levels are considered external to the model. We therefore focus on the number of individuals as opposed to the number of genome copies. Sizes and mating system details are specified for each population within population-specific epochs. Epochs are contiguous time intervals that define the existence interval of the population. Each epoch specifies the population size over that interval, which can be a constant value or a function defined by start and end sizes that must remain positive.

A1.3 Population dynamics

Within a population, we assume that allele frequency dynamics can be described by the Wright-Fisher model. Briefly, generations are non-overlapping (all parents reproduce and die simultaneously), and for allele i currently at frequency p_i , its frequency in the next generation (at birth) is expected to be $p_i w_i / \bar{w}$, where w_i and \bar{w} are the marginal and mean fitnesses, respectively, properly weighted according to ancestry proportions. In this framework, a forward-time simulation of finite populations is equivalent to multinomial sampling of allele frequencies each generation (BÜRGER (2000, pp 29-31), CROW and KIMURA (1970, pp 179-181)), and a backwards-time (coalescent) simulation follows the approximations described in TAJIMA (1983), HUDSON (1983) and WAKELEY (2008, chapter 3). Further, this model assumes “soft” selection (CHRISTIANSEN, 1975), meaning that the dynamics of population sizes changes are independent of the details of individual fitnesses. As such, this model excludes scenarios such as “hard selection,” in which population sizes are dependent on a population’s mean fitness, or stochastic fluctuations in population size, such as interpreting population sizes as carrying capacities. Many forwards and backwards time simulators currently implement this model (e.g., HUDSON, 2002, GUTENKUNST *et al.*, 2009, EXCOFFIER and FOLL, 2011, KELLEHER *et al.*, 2016, JOUGANOUS *et al.*, 2017, HALLER and MESSER, 2019, THORNTON, 2019). Each has their own unique method of specifying the details of selection and demography. The goal of Demes is to provide a common method for the latter task.

A1.4 Selfing and cloning

Each population has an assigned selfing rate and cloning rate, where each defines the probability that offspring are generated from one generation to the next by either self-fertilization or cloning of an individual. More specifically, for a given epoch within a population denote the clonal rate by σ and the selfing rate by S . S and σ can take any value between zero and one and can sum to more than one. Each generation a proportion of offspring σ are expected to be generated through clonal reproduction, while $1 - \sigma$ are expected to arise through sexual reproduction. Within the sexually-reproduced offspring, S are born via self-fertilization while the rest have parents drawn at random from the previous generation. Depending on the simulator, this random drawing of parent may occur either with or without replacement. When drawing occurs with replacement, a small amount of “residual” selfing is expected, so that the realized selfing probability is $(1 - \sigma)(S +$

$(1 - S)/N$) instead of $(1 - \sigma)S$ (so that even with $\sigma = 0$ and $S = 0$, selfing may still occur with probability $1/N$), although this effect is negligible in large populations (NORDBORG and DONNELLY, 1997). Simulators that allow variable rates of selfing are expected to clearly document the expected behavior.

By allowing the definition of selfing and cloning probabilities, we allow many standard models to be defined. However, by parameterizing selfing and cloning as we have, we assume that these properties of populations can be specified independently from the genetics. In other words, mutations that cause selfing probabilities to fluctuate within an epoch are not considered. More details of the mathematical properties of selfing and cloning rates in a coalescent context can be found in NORDBORG and DONNELLY (1997), HARTFIELD *et al.* (2016).

A1.5 Relationships between populations

A population may have one or more ancestors, which are other populations that exist at the population's start time. If one ancestor is specified, the first generation is constructed by randomly sampling parents from the ancestral population to contribute to offspring in the newly generated population. If more than one ancestor is specified, the proportions of ancestry from each contributing population must be provided, and those proportions must sum to one. In this case, parents are chosen randomly from each ancestral population with probability given by those proportions.

Individuals in a population may have parents from a different population through migrations. These can be defined as continuous migration rates over time intervals for which populations co-exist or through instantaneous (or pulse) migration events at a given time. Continuous migration rates are defined as the probability that parents in the "destination" population are chosen from the "source" population. On the other hand, pulse migration events specify the instantaneous replacement of a given fraction of individuals in a destination population by individuals with parents from a source population.

A2 Human Data Model

The Demes Human Data Model (HDM) is focused on the efficiency of human understanding and the avoidance of errors. We have adopted the widely used YAML format (BEN-KIKI *et al.*, 2009) as the recommended means of interchanging Demes models (e.g., Figures A1 and A2). YAML is a data serialization language with an emphasis on simplicity and which interoperates well with JSON (indeed, YAML 1.2 is a super set of JSON). We chose YAML over JSON because although JSON is an excellent format for data interchange, it is ill-suited for human understanding or manipulation. We also considered other declarative data exchange formats such as TOML, but chose YAML because of equivalence with JSON, its popularity and good software support. Since the Demes data model is defined in JSON Schema, however, there is no formal dependency on YAML and implementations may choose to use JSON directly if they wish (e.g., for greater efficiency).

Structurally, the HDM encourages human understanding by avoiding redundancy in the description where possible and by providing a mechanism for specifying default values that are inherited hierarchically. Figure A1 shows an example Demes model expressed in both the HDM and Machine Data Model (MDM) forms (both in terms of YAML syntax). For values that repeat across fields, defaults may be used to implicitly assign default values to fields of the given type. A default is superseded by an explicitly provided value if given. Other implicit values are inherited naturally

following the progression of time. For example, if an epoch `start_size` is not provided, it is assumed to be equal to the `end_size` of the previous epoch. This also means that the first epoch of each population must specify the initial size.

A3 Parsers

The implicit value inference and default value propagation required to fully resolve a Demes HDM description into the corresponding MDM form is straightforward to describe, but still requires significant effort to implement in software. Moreover, there are many constraints on the data model (for example, epoch `end_time` values should be strictly decreasing within a deme), which must be enforced. It would not be reasonable to require every program that takes Demes models as input to implement this logic, as the programming effort would be significant (limiting adoption) and it is likely that if there were many implementations some would differ in their details (harming the software ecosystem). We therefore define a standard software entity as part of the specification (the parser), which performs this task and can be shared by programs that support Demes as an input. By having relatively few high-quality Demes parsers available as libraries (ideally, one per programming language), the probability of divergences from the standard is greatly reduced.

The Demes specification precisely defines the required behavior of parsers, which translate a model described in the HDM into the MDM, where values have been assigned to all parameters and data constraints have been checked. The output is formally defined as JSON, but in practice parsers output an object model that is suitable for the particular programming language and that can be used directly by the implementing program. We provide a reference implementation written in Python to resolve any potential ambiguities and to provide a helpful template for other implementations, as well as an extensive test suite of examples and the expected outputs. In addition, we have high-quality parser implementations in the Python and C languages (published under liberal open source licenses), providing a solid foundation for the software ecosystem.

A4 Scope of the specification

A4.1 Static models, not parameterized models

The Demes specification is designed to describe demographic models defined by a fixed set of model parameters. As described in the main text, it does not include information about estimated confidence intervals or the joint distribution of parameter values. In this section we describe the rationale for this design decision.

The parameters of demographic models are typically tightly coupled, and cases in which distributions for different parameters can be simply described are rare. In this situation, the simplest way to describe an estimated distribution is to list a large number of samples from the posterior. While writing out a large number of Demes models in YAML format may seem inefficient, it can in fact be a compact way to describe these distributions. For example, consider a one-population model with piecewise-constant sizes over 20 epochs which has ~ 40 free parameters: the `start_size` and `end_time` values for each epoch. If we sample 50,000 models from the posterior distribution, the resulting multi-document YAML file is 45 MiB. This format compresses down to 8.4 MiB when gzipped or 6.2 MiB when compressed with LZMA2, which is on par with an equivalent binary representation of the free parameters ($40 \times 50000 \times 4 \text{ bytes} \approx 7.6 \text{ MiB}$).

Similarly, one might be interested in running simulations in which the demographic model parameters are drawn from a distribution, e.g., in ABC inference (BEAUMONT *et al.*, 2002). Other inference procedures based on optimizing a loss function (GUTENKUNST *et al.*, 2009, KAMM *et al.*, 2017, JOUGANOUS *et al.*, 2017, RAGSDALE and GRAVEL, 2019, EXCOFFIER *et al.*, 2021) need users to specify parameter bounds, and possibly non-linear or conditional constraints between parameters. Indeed, the choice of how to parameterize a model could be important for some inference methods (e.g. absolute times versus relative times between events).

Implementing the many distributions of interest and supporting a general way to describe a model's free parameters would greatly increase the complexity of parsers, with relatively limited benefit to most users. It's unlikely that Demes could be made sufficiently flexible without considerable effort to implement many features of general-purpose programming languages, such as variables, arithmetic, and flow control. Such use cases are therefore better served by writing model-generating functions in an existing programming language, for example using the Demes Python API (e.g., as implemented in *moments* (JOUGANOUS *et al.*, 2017, RAGSDALE and GRAVEL, 2019)). As an intriguing possibility for developments in this direction, there exist many templating solutions for YAML and JSON that are specifically designed for extending static data in arbitrarily complex ways (e.g., YTT, Jsonnet, CUE, and Dhall).

A4.2 Population-level features, not genome features

Demes is designed to interface with a large number of simulation and inference methods, so we have restricted the specification to only describe demographic features at the population level. Features of the underlying genome biology are omitted, including mutation and recombination rates, genome annotations, ploidy, sequence annotations, and so on. Selection and dominance models are absent, as discussed in Sec. A1, and we do not specify how individuals are sampled.

The model is, by design, very limited. We do not specify details of individuals or of the processes that happen within populations, because such details are necessarily complicated, vary greatly by application, and attempts to encompass such a broad range of biological processes are unlikely to succeed. The narrow focus on individuals and their groupings into populations should ensure that the Demes standard is stable, and not subject to continual development as more and more elaborations of basic processes are added. Demes is intended to provide *part* of a simulation model specification; other parts of the specification, such as sampling strategies, or parameters that vary by time or population can *refer to* this well defined and unambiguous description of the demography.

A5 Extended data and figures

C

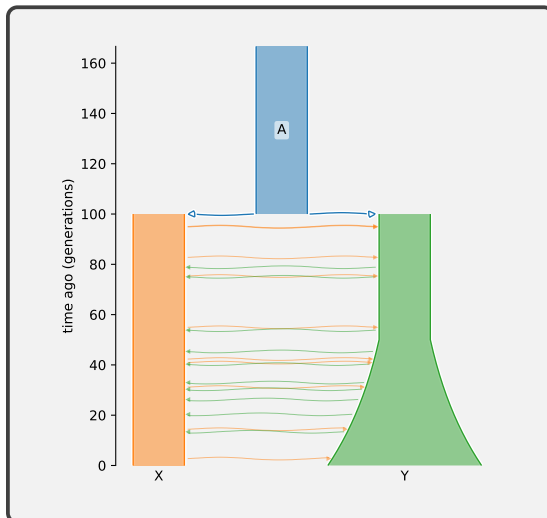
A

```

1 # Comments start with a hash.
2 description:
3   Two-deme isolation-with-migration model.
4 time_units: generations
5 defaults:
6   epoch: {start_size: 1000}
7 demes:
8   - name: A
9     description: The ancestral deme
10    epochs:
11      - end_time: 100
12   - name: X
13     description: First descendant deme.
14     ancestors: [A]
15   - name: Y
16     description: Second descendant deme.
17     ancestors: [A]
18     epochs:
19       - {end_time: 50}
20       - {end_size: 3000}
21 migrations:
22   - {demes: [X, Y], rate: 1e-4}

```

B



```

1 description: Two-deme isolation-with-migration model.
2 time_units: generations
3 generation_time: 1
4 doi: []
5 demes:
6 - name: A
7   description: The ancestral deme
8   start_time: .inf
9   ancestors: []
10  proportions: []
11  epochs:
12 - end_time: 100
13   start_size: 1000
14   end_size: 1000
15   size_function: constant
16   selfing_rate: 0
17   cloning_rate: 0
18 - name: X
19   description: First descendant deme.
20   start_time: 100
21   ancestors: [A]
22   proportions: [1]
23   epochs:
24 - end_time: 0
25   start_size: 1000
26   end_size: 1000
27   size_function: constant
28   selfing_rate: 0
29   cloning_rate: 0
30 - name: Y
31   description: Second descendant deme.
32   start_time: 100
33   ancestors: [A]
34   proportions: [1]
35   epochs:
36 - end_time: 50
37   start_size: 1000
38   end_size: 1000
39   size_function: constant
40   selfing_rate: 0
41   cloning_rate: 0
42 - end_time: 0
43   start_size: 1000
44   end_size: 3000
45   size_function: exponential
46   selfing_rate: 0
47   cloning_rate: 0
48 migrations:
49 - {source: X, dest: Y, start_time: 100, end_time: 0,
50   rate: 0.0001}
51 - {source: Y, dest: X, start_time: 100, end_time: 0,
52   rate: 0.0001}
53 pulses: []

```

Figure A1: Example isolation-with-migration Demes model. (A) The Human Data Model representation expressed using YAML. (B) A visual representation of the model using demesdraw. (C) The same model in the Machine Data Model form. The YAML descriptions in (A) and (C) correspond exactly to a JSON description, but are much more human-readable.


```
1 description: The two-population model inferred in Tennesen et al (2012).
2 doi: ["https://doi.org/10.1038/nature11690"]
3 time_units: years
4 generation_time: 25
5 demes:
6 - name: ancestral
7   description: Population that splits into EUR and AFR
8   epochs:
9     - start_size: 7310
10     end_time: 148000
11     - start_size: 14474
12     end_time: 51000
13 - name: AFR
14   description: African Americans
15   ancestors: [ancestral]
16   epochs:
17     - start_size: 14474
18     end_time: 5115
19     - end_time: 0
20     end_size: 432125
21 - name: EUR
22   description: European Americans
23   ancestors: [ancestral]
24   epochs:
25     - start_size: 1861
26     end_time: 23000
27     - start_size: 1032
28     end_time: 5115
29     end_size: 9279
30     - end_time: 0
31     end_size: 501436
32 migrations:
33 - demes: [AFR, EUR]
34   rate: 1.5e-4
35   end_time: 5115
36 - demes: [AFR, EUR]
37   rate: 2.5e-5
38   start_time: 5115
```

Figure A2: **The TENNESSEN *et al.* (2012) two-population demographic model in Demes format.** This model includes a single ancestral population that expands in size in the past, followed by divergence between AFR- and EUR-labeled populations. The two-population phase of the model includes multiple epochs of varying size, and rapid exponential growth over the past five thousand years in each population.

```
1 import demes
2 import msprime
3 import moments
4
5 #####
6 ##### Import the demographic model using demes, set up samples
7 #####
8
9 graph = demes.load("tennessen.yaml")
10 samples = {"AFR": 20, "EUR": 20}
11
12 #####
13 ##### Simulate genomic data using msprime
14 #####
15
16 msprime_samples = [
17     msprime.SampleSet(n, ploidy=1, population=p) for p, n in samples.items()
18 ]
19 demog = msprime.Demography.from_demes(graph) # load the demes model as msprime demography
20 L = 1e7 # sequence length of 10 Mb
21 r = 2e-8 # with constant recombination rate of 2e-8
22 u = 2.36e-8 # mutation rate used in Tennessen et al (2012)
23
24 ts = msprime.sim_ancestry(
25     msprime_samples,
26     demography=demog,
27     sequence_length=L,
28     recombination_rate=r,
29     random_seed=1234567,
30 )
31
32 ts = msprime.sim_mutations(ts, rate=u, random_seed=1234567)
33
34 # compute the SFS from the msprime simulation using tskit
35 msprime_afr = ts.allele_frequency_spectrum(
36     [range(samples["AFR"])], mode="site", polarised=True, span_normalise=False
37 )
38 msprime_eur = ts.allele_frequency_spectrum(
39     [range(samples["AFR"], samples["AFR"] + samples["EUR"])],
40     mode="site",
41     polarised=True,
42     span_normalise=False,
43 )
44
45 #####
46 ##### Compute expected SFS for sampled populations using moments
47 #####
48
49 Ne = graph["ancestral"].epochs[0].start_size
50 theta = 4 * Ne * u * L
51
52 fs = moments.Spectrum.from_demes(graph, samples=samples)
53 fs *= theta # rescale to match the total mutation rate in the msprime simulation
54
55 moments_afr = fs.marginalize([1])
56 moments_eur = fs.marginalize([0])
```

Figure A3: **Simulation of SFS for the Tennessen model.** We first load the demographic model using `demes` (as `graph`), which can then be used by `msprime` to create the demographic model used in `msprime.sim_ancestry()`. The same loaded graph can also be passed to `moments` to compute the expected joint SFS. To compare the SFS in Figure 1, we *marginalize* the joint SFS to obtain the single-population SFS for both AFR and EUR populations. Lines interfacing `demes` and other software are highlighted.