

1 Running title: ipcoal: simulate genealogies and sequences on species trees

2

## 3 **ipcoal: An interactive Python package for simulating and analyzing** 4 **genealogies and sequences on a species tree or network**

5 Patrick F. McKenzie<sup>1,2</sup> and Deren A. R. Eaton<sup>1</sup>

6 <sup>1</sup> Department of Ecology, Evolution, and Environmental Biology, Columbia University, New York, NY  
7 10027

8 <sup>2</sup> To whom correspondence should be addressed

### 9 **Abstract–**

10 **Summary:** *ipcoal* is a free and open source Python package for simulating and analyzing genealogies and  
11 sequences. It automates the task of describing complex demographic models (e.g., with divergence times,  
12 effective population sizes, migration events) to the *msprime* coalescent simulator by parsing a user-supplied  
13 species tree or network. Genealogies, sequences, and metadata are returned in tabular format allowing for  
14 easy downstream analyses. *ipcoal* includes phylogenetic inference tools to automate gene tree inference  
15 from simulated sequence data, and visualization tools for analyzing results and verifying model accuracy.  
16 The *ipcoal* package is a powerful tool for posterior predictive data analysis, for methods validation, and for  
17 teaching coalescent methods in an interactive and visual environment.

18  
19 **Availability and implementation:** Source code is available from the GitHub repository ([https://github.com/  
20 pmckenz1/ipcoal/](https://github.com/pmckenz1/ipcoal/)) and is distributed for packaged installation with conda. Complete documentation and in-  
21 teractive notebooks prepared for teaching purposes are available at <https://ipcoal.readthedocs.io/>.

22  
23 **Keywords:** coalescent, evolution, simulation, Python, phylogeny

## 24 **1 Introduction**

25 The coalescent process (Hudson, 1983; Kingman, 1982) is used to model the distribution of genealogi-  
26 cal ancestry across a set of sampled genomes. It approximates a neutral Wright-Fisher process of random  
27 mating within populations where the expected waiting times between subsequent coalescent events can be  
28 drawn from a statistical distribution based on the effective population size. This makes simulation of ge-  
29 nealogies under the coalescent process (Hudson, 2002) a computationally efficient approach for integrating  
30 over genealogical variation (i.e., treating it as a latent random variable) when making population genetic  
31 inferences (Beerli & Felsenstein, 2001).

32 Demographic models specify the parameters of a coalescent simulation. Highly complex models may  
33 include population sizes and divergence times, and gene flow (admixture) between populations. For ex-  
34 ample, in the study of human history, a demographic model may describe divergences among different  
35 continents, the expansion of populations separately in Africa, Eurasia, and the Americas, and subsequent  
36 admixture between them (Reich, 2018; Gronau *et al.*, 2011; Green *et al.*, 2010). Demographic models are

37 also routinely used in phylogenetics, with the goal of inferring a topology (i.e., the relationships among  
38 connected populations) in addition to the parameters of a demographic model applied to the topology  
39 (Knowles & Kubatko, 2011; Degnan & Rosenberg, 2009).

40 The ability to simulate realistic sequence data evolving on genealogies sampled from complex demo-  
41 graphic models has enabled new types of inference from genomic data, from fitting parameters to demo-  
42 graphic models and performing model comparisons (Chung & Hey, 2017); to performing posterior predic-  
43 tive data analyses (Brown, 2014); to generating training datasets for machine learning methods (Schri-  
44 der & Kern, 2018); to validating new inference methods (Adrion *et al.*, 2019). Despite the impressive capa-  
45 bilities of recent state-of-the-art coalescent simulation tools like *msprime* (Kelleher *et al.*, 2016), it is dif-  
46 ficult for a single package to be optimized for all types of use. To this end, *msprime* lacks functionality in  
47 ways that limit its utility for studying deeper-scale (e.g., phylogenetic) datasets. Here we describe a new  
48 Python package, *ipcoal*, which wraps around *msprime* with the aim of filling this niche: to provide a sim-  
49 ple method for simulating genealogies and sequences on species trees or networks.

## 50 2 Phylogenomic data simulation

51 We make the following distinctions among terms in *ipcoal*: a genealogy is the true history of ancestry  
52 among a set of sampled genes; a gene tree is an empirical estimate of a genealogy based on sequences  
53 from some region of the genome; and a species tree is a demographic model including a topology (Mad-  
54 dison, 1997; Pamilo & Nei, 1988). As phylogenetics transitions from a focus on multi-locus data sets  
55 (Knowles & Kubatko, 2011) to the analysis of whole genomes – and the spatial distribution of correlated  
56 genealogical variation along chromosomes – these distinctions that we highlight in *ipcoal*, between un-  
57 observable genealogical variation and the empirical gene tree estimates that can be made from observable  
58 sequence data, will become increasingly relevant (Adams & Castoe, 2019; Posada & Crandall, 2002).

59 Simulating realistic sequence data under the multispecies coalescent model has typically involved a  
60 two-step approach: a set of independent genealogies is first simulated, and then a model of sequence evo-  
61 lution is applied along the edges of each tree to produce sequence alignments. This phylogenetic workflow  
62 differs from standard population-level coalescent simulations in several ways: (1) phylogenies generally  
63 contain many more lineages than population genetic models which makes describing them to coalescent  
64 simulators burdensome and error-prone; (2) the phylogenetic workflow typically ignores recombination,  
65 but such data can now be simulated easily by modern coalescent software; and (3) the phylogenetic work-  
66 flow applies a Markov model of sequence evolution rather than the more simple infinite-sites process, al-  
67 lowing for homoplasy and asymmetrical substitution rates. In *ipcoal* we have combined the best aspects of  
68 each approach so that it is easy to describe demographic models for large trees, to simulate independent or  
69 linked genealogies, and to generate sequences under complex models of sequence evolution.

## 70 3 Implementation

### 71 3.1 Reproducible and robust workflow

72 The *ipcoal* library is designed for interactive use within jupyter-notebooks (Kluyver *et al.*, 2016), where  
73 simulations can be run in the same document as downstream statistical analyses; visualization tools can be  
74 used to validate model accuracy; and code, figures, and results are easily organized into reproducible and

75 shareable documents. The code is designed to be easy to use, following a minimalist and object-oriented  
76 design with few user-facing classes and functions.

## 77 **3.2 Defining demographic models**

78 The primary object that users interact with in *ipcoal* is the `Model` class object (Fig. 1a), which takes a  
79 number of user-supplied parameter arguments to initialize demographic and substitution models. The pri-  
80 mary convenience of the `Model` class object is its ability to automate the construction of a demographic  
81 model by parsing a tree object. For large phylogenies this is important. For example, to describe a de-  
82 mographic model for a species tree with 20 tips in *msprime* would require writing code to define 39 di-  
83 vergence events (`MassMigrations`). *ipcoal* uses the Python tree manipulation and plotting library *toytree*  
84 (Eaton, 2020) to parse, visualize, and annotate trees, making it easy to verify whether variable  $N_e$  values  
85 and admixture scenarios have been properly defined (Fig. 1a-b).

## 86 **3.3 Simulating unlinked SNPs**

87 Many inference tools require the input of unlinked single nucleotide polymorphisms (SNPs) to circumvent  
88 the effect of recombination (e.g., SVDquartets (Chifman & Kubatko, 2014) and SNAPP (Bryant *et al.*,  
89 2012)). *ipcoal* can generate a distribution of independent genealogies, and unlinked SNPs evolved on those  
90 genealogies, using the `Model.sim_snps()` function call (Fig. 1c-d). Notably, we take care that the  
91 probability with which a substitution is observed is proportional to the total edge lengths of the genealogy  
92 by testing each genealogy for a SNP and moving on to the next independently sampled genealogy if a SNP  
93 is not observed. By contrast, users can alternatively toggle the option to enforce a SNP placement on ev-  
94 ery visited genealogy, which will increase the speed of simulations but introduce a bias toward shallower  
95 divergence times.

## 96 **3.4 Simulating loci**

97 The `Model` object can also simulate entire chromosomes (loci) with or without recombination by calling  
98 the `Model.sim_loci()` function. This produces sequences of linked genealogies. Nearby genealogies  
99 are correlated since some samples share the same ancestors at neighboring genomic regions, and thus are  
100 more similar in topology and edge lengths than unlinked trees (Fig. 1d). This type of variation is increas-  
101 ingly of interest for genome-wide analyses.

## 102 **3.5 Simulating sequence evolution**

103 To simulate sequence data on genealogies in *ipcoal*, a continuous-time Markov substitution model is ap-  
104 plied iteratively to each edge of the tree from root to tips. We have implemented our own sequence simu-  
105 lator using just-in-time compiled Python code to achieve high performance. We additionally provide the  
106 option of using the external tool *seq-sen* (Rambaut & Grass, 1997), which offers a larger range of models  
107 than we currently support. Our internal implementation is used by default since it achieves faster speeds  
108 by avoiding repeated subprocess calls. The documentation includes test notebooks demonstrating that our  
109 implementation converges to the same results as *seq-sen*.

## 110 3.6 Results

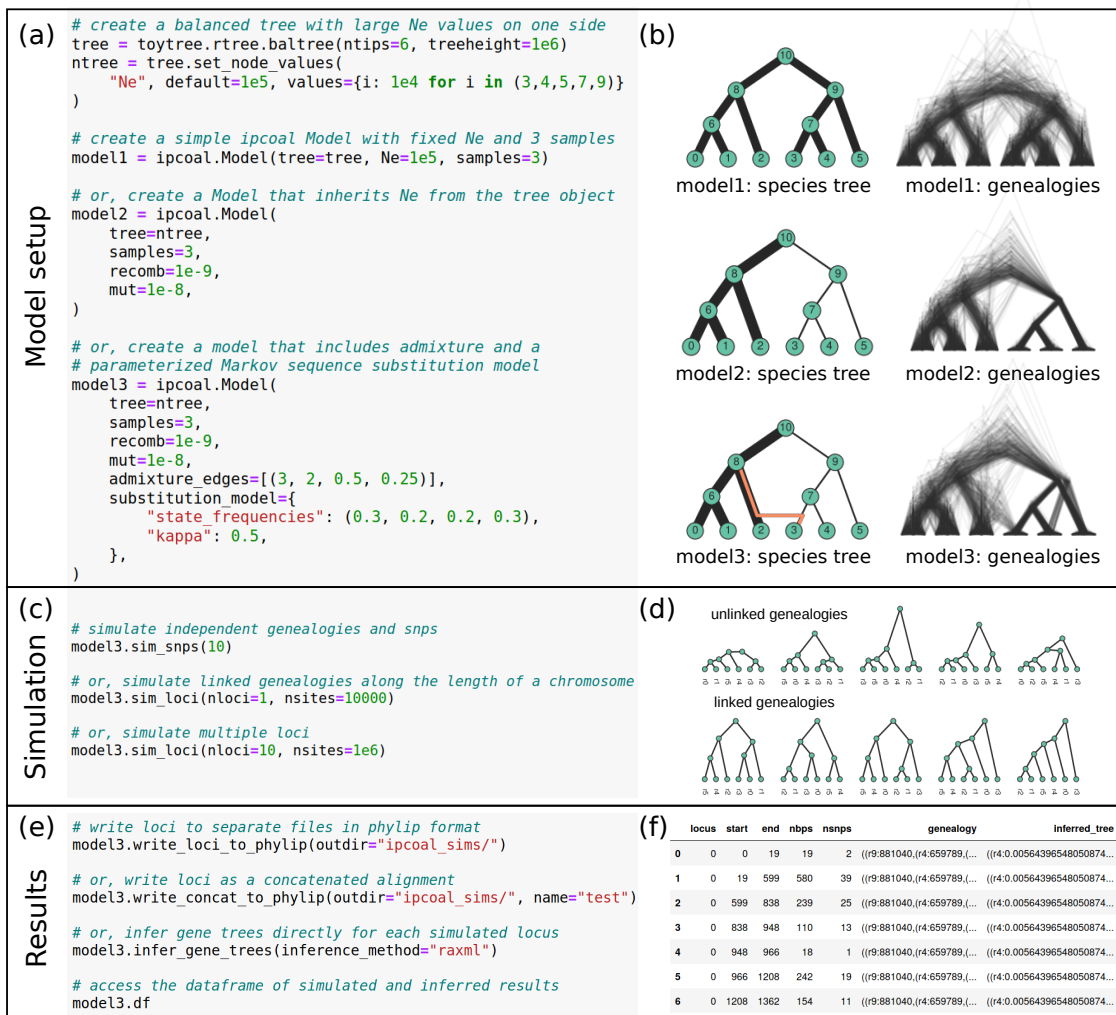
111 Upon calling a simulation function, two results are stored to the `Model` object: a sequence array (`Model.seqs`)  
112 and a dataframe with the genealogy and statistics about each genealogical window (`Model.df`). The se-  
113 quence array can be written to disk in Nexus or Phylip format, and as separate or concatenated loci, and  
114 the DataFrame can be saved as a CSV (Fig. 1e-f). However, to simplify analytical workflows, we provide  
115 convenience functions for inferring gene trees directly from sequence data, avoiding the need to organize  
116 many files.

## 117 4 Conclusions

118 Coalescent simulations for studying genome-wide patterns are routinely used in population genetics, but  
119 have not yet achieved widespread use in phylogenetics where the focus has traditionally been limited to  
120 a smaller number of unlinked loci. Our new software tool *ipcoal* makes it easy to simulate and explore  
121 linked or unlinked genealogical and sequence variation across genomes, providing new opportunities for  
122 investigating phylogenetic methods and theory.

## 123 References

- 124 Adams, R.H. & Castoe, T.A. (2019). Statistical binning leads to profound model violation due to gene tree  
125 error incurred by trying to avoid gene tree error. *Molecular Phylogenetics and Evolution*, 134, 164 –  
126 171. [2](#)
- 127 Adrion, J.R., Cole, C.B., Dukler, N., Galloway, J.G., Gladstein, A.L., Gower, G., Kyriazis, C.C., Ragsdale,  
128 A.P., Tsambos, G., Baumdicker, F., Carlson, J., Cartwright, R.A., Durvasula, A., Kim, B.Y., McKen-  
129 zie, P., Messer, P.W., Noskova, E., Vecchy, D.O.D., Racimo, F., Struck, T.J., Gravel, S., Gutenkunst,  
130 R.N., Lohmeuller, K.E., Ralph, P.L., Schrider, D.R., Siepel, A., Kelleher, J. & Kern, A.D. (2019). A  
131 community-maintained standard library of population genetic models. *bioRxiv*, p. 2019.12.20.885129. [2](#)
- 132 Beerli, P. & Felsenstein, J. (2001). Maximum likelihood estimation of a migration matrix and effective  
133 population sizes in n subpopulations by using a coalescent approach. *Proceedings of the National*  
134 *Academy of Sciences*, 98, 4563–4568. [1](#)
- 135 Brown, J.M. (2014). Predictive Approaches to Assessing the Fit of Evolutionary Models. *Systematic*  
136 *Biology*, 63, 289–292. [2](#)
- 137 Bryant, D., Bouckaert, R., Felsenstein, J., Rosenberg, N.A. & RoyChoudhury, A. (2012). Inferring Species  
138 Trees Directly from Biallelic Genetic Markers: Bypassing Gene Trees in a Full Coalescent Analysis.  
139 *Molecular Biology and Evolution*, 29, 1917–1932. [3](#)
- 140 Chifman, J. & Kubatko, L. (2014). Quartet Inference from SNP Data Under the Coalescent Model. *Bioin-*  
141 *formatics*, 30, 3317–3324. [3](#)
- 142 Chung, Y. & Hey, J. (2017). Bayesian Analysis of Evolutionary Divergence with Genomic Data under  
143 Diverse Demographic Models. *Molecular Biology and Evolution*, 34, 1517–1528. [2](#)



**Figure 1.** Simulation of coalescent genealogies and sequence data in *ipcoal*. A species tree can be generated or loaded from a newick string to define population relationships in a demographic model, and a single  $N_e$  value or variable  $N_e$ s can be applied to nodes by mapping values using *toytree*. The Model class object of *ipcoal* is used to initialize parameterized demographic and mutational models (a). Genealogical variation reflects parameters of the demographic model including  $N_e$  and admixture events, each of which can be easily visualized for validation (b). Sequence data can be simulated as unlinked SNPs (c) or as continuous loci in which recombination affects linkage among neighboring genealogies (d). Simulated sequences can be written to files (either concatenated or as separate loci) for downstream analyses, or the sequences can be used to directly infer gene trees (e). Simulated and inferred results are organized into dataframes for further analyses (f).

- 144 Degnan, J.H. & Rosenberg, N.A. (2009). Gene tree discordance, phylogenetic inference and the multi-  
145 species coalescent. *Trends in Ecology & Evolution*, 24, 332–340. 2
- 146 Eaton, D.A.R. (2020). Toytree: A minimalist tree visualization and manipulation library for Python. *Meth-*  
147 *ods in Ecology and Evolution*, 11, 187–191. 3
- 148 Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li, H., Zhai, W.,  
149 Fritz, M.H.Y., Hansen, N.F., Durand, E.Y., Malaspinas, A.S., Jensen, J.D., Marques-Bonet, T., Alkan,  
150 C., Prüfer, K., Meyer, M., Burbano, H.A., Good, J.M., Schultz, R., Aximu-Petri, A., Butthof, A., Höber,  
151 B., Höffner, B., Siegemund, M., Weihmann, A., Nusbaum, C., Lander, E.S., Russ, C., Novod, N., Af-  
152 fourtit, J., Egholm, M., Verna, C., Rudan, P., Brajkovic, D., Kucan, Ž., Gušić, I., Doronichev, V.B.,  
153 Golovanova, L.V., Lalueza-Fox, C., de la Rasilla, M., Fortea, J., Rosas, A., Schmitz, R.W., Johnson,  
154 P.L.F., Eichler, E.E., Falush, D., Birney, E., Mullikin, J.C., Slatkin, M., Nielsen, R., Kelso, J., Lachmann,  
155 M., Reich, D. & Pääbo, S. (2010). A draft sequence of the neandertal genome. *Science*, 328, 710–722. 1
- 156 Gronau, I., Hubisz, M.J., Gulko, B., Danko, C.G. & Siepel, A. (2011). Bayesian inference of ancient hu-  
157 man demography from individual genome sequences. *Nature Genetics*, 43, 1031–1035. 1
- 158 Hudson, R.R. (1983). Testing the Constant-Rate Neutral Allele Model with Protein Sequence Data. *Evolu-*  
159 *tion*, 37, 203–217. 1
- 160 Hudson, R.R. (2002). Generating samples under a Wright–Fisher neutral model of genetic variation.  
161 *Bioinformatics*, 18, 337–338. 1
- 162 Kelleher, J., Etheridge, A.M. & McVean, G. (2016). Efficient Coalescent Simulation and Genealogical  
163 Analysis for Large Sample Sizes. *PLOS Computational Biology*, 12, e1004842. 2
- 164 Kingman, J.F.C. (1982). The coalescent. *Stochastic Processes and their Applications*, 13, 235–248. 1
- 165 Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B.E., Bussonnier, M., Frederic, J., Kelley, K., Hamrick,  
166 J.B., Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., Willing, C. & al, e. (2016). Jupyter Note-  
167 books - a publishing format for reproducible computational workflows. In: *ELPUB*. 2
- 168 Knowles, L.L. & Kubatko, L.S. (eds.) (2011). *Estimating Species Trees: Practical and Theoretical As-*  
169 *pects*. 1st edn. Wiley-Blackwell. 2
- 170 Maddison, W.P. (1997). Gene Trees in Species Trees. *Systematic Biology*, 46, 523–536. 2
- 171 Pamilo, P. & Nei, M. (1988). Relationships between gene trees and species trees. *Molecular Biology and*  
172 *Evolution*, 5, 568–583. 2
- 173 Posada, D. & Crandall, K.A. (2002). The effect of recombination on the accuracy of phylogeny estimation.  
174 *Journal of Molecular Evolution*, 54, 396–402. 2
- 175 Rambaut, A. & Grass, N.C. (1997). Seq-Gen: an application for the Monte Carlo simulation of DNA se-  
176 quence evolution along phylogenetic trees. *Bioinformatics*, 13, 235–238. 3
- 177 Reich, D. (2018). *Who we are and how we got here: Ancient DNA and the new science of the human past*.  
178 Oxford University Press. 1

179 Schrider, D.R. & Kern, A.D. (2018). Supervised machine learning for population genetics: A new  
180 paradigm. *Trends in Genetics*, 34, 301 – 312. [2](#)