OXFORD GENETICS

Advance Access Publication Date: 6 March 2025 Investigation

Evaluating ARG-estimation methods in the context of estimating population-mean polygenic score histories

Dandan Peng, Obadiah J. Mulder, Michael D. Edge 🕩 *

Department of Quantitative and Computational Biology, University of Southern California, 1050 Childs Way, Los Angeles, CA 90098, USA

*Corresponding author: Department of Quantitative and Computational Biology, University of Southern California, 1050 Childs Way, Los Angeles, CA 90098, USA. Email: edgem@usc.edu

Scalable methods for estimating marginal coalescent trees across the genome present new opportunities for studying evolution and have generated considerable excitement, with new methods extending scalability to thousands of samples. Benchmarking of the available methods has revealed general tradeoffs between accuracy and scalability, but performance in downstream applications has not always been easily predictable from general performance measures, suggesting that specific features of the ancestral recombination graph (ARG) may be important for specific downstream applications of estimated ARGs. To exemplify this point, we benchmark ARG estimation methods with respect to a specific set of methods for estimating the historical time course of a population-mean polygenic score (PGS) using the marginal coalescent trees encoded by the ARG. Here, we examine the performance in simulation of seven ARG estimation methods: ARGweaver, RENT+, Relate, tsinfer+tsdate, ARG-Needle, ASMC-clust, and SINGER, using their estimated coalescent trees and examining bias, mean squared error, confidence interval coverage, and Type I and II error rates of the down-stream methods. Although it does not scale to the sample sizes attainable by other new methods, SINGER produced the most accurate estimated PGS histories in many instances, even when Relate, tsinfer+tsdate, ARG-Needle, and ASMC-clust used samples 10 or more times as large as those used by SINGER. In general, the best choice of method depends on the number of samples available and the historical time period of interest. In particular, the unprecedented sample sizes allowed by Relate, tsinfer+tsdate, ARG-Needle, and ASMC-clust are of greatest importance when the recent past is of interest—further back in time, most of the tree has coalesced, and differences in contemporary sample size are less salient.

Keywords: ancestral recombination graphs; coalescent; polygenic traits; natural selection

Introduction

The ancestral recombination graph, or ARG (Griffiths and Marjoram 1996), is a rich representation of the history of a sample of haplotypes, including all the mutation, recombination, and common-ancestry events that affect contemporary variation. Thus, the ARG encodes all historical information that can be extracted from a sample of contemporary genomes, much of it in gene genealogies or coalescent trees (Hudson 1990; Wakeley 2016) for every location in the genome, termed "local" or "marginal" trees. The true ARG is generally unknown, and estimation of the ARG is a very challenging problem. Nonetheless, the last 10 years have witnessed major advances in ARG estimation, with new methods that provide estimated ARGs with unprecedented accuracy, scalability, or both (Rasmussen et al. 2014; Mirzaei and Wu 2017; Kelleher et al. 2019; Speidel et al. 2019; Zhang et al. 2023; Deng et al. 2024). These advances have produced a great deal of excitement about the potential of estimated ARGs in evolutionary biology and beyond (Harris 2019, 2023; Brandt et al. 2024; Lewanski et al. 2024; Wong et al. 2024; Nielsen et al. 2025).

The promise of estimated ARGs depends on their performance in downstream applications. Many of the available $\ensuremath{\mathsf{ARG}}$

estimation methods have been benchmarked in general terms (Deng *et al.* 2021; Brandt *et al.* 2022), revealing tradeoffs between accuracy and scalability. (We refer to all the methods we consider here as "ARG estimation" methods, regardless of whether they estimate the full ARG, including recombinations between marginal trees, or just a set of marginal trees.) However, early indications are that the actual performance of ARG estimators in downstream applications can vary in ways that are not necessarily predicted by a generic accuracy-versus-scalability tradeoff (Fan *et al.* 2022, 2023). Thus, it seems as though the performance of estimated ARGs in downstream procedures may depend on specific features of the ARG and how well they are estimated.

To explore this point in depth, we conducted thorough benchmarking of ARG estimators with respect to a set of methods for studying polygenic traits. These methods take estimated marginal trees as input. Polygenic traits—traits influenced by genetic variants from across the genome—are a promising area for applications of estimated ARGs (Edge and Coop 2019; Stern *et al.* 2021; Link *et al.* 2023; Zhang *et al.* 2023; Christ *et al.* 2024; Gunnarsson *et al.* 2024; Wang *et al.* 2024; Zhu *et al.* 2024). Understanding the evolution of polygenic traits is challenging in part because signals of selection are spread across many loci. Researchers can study

© The Author(s) 2025. Published by Oxford University Press on behalf of The Genetics Society of America.

Received on 07 January 2025; accepted on 15 February 2025

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs licence (https://creativecommons. org/licenses/by-nc-nd/4.0/), which permits non-commercial reproduction and distribution of the work, in any medium, provided the original work is not altered or transformed in any way, and that the work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

the history of polygenic traits by examining fossil records (Kappelman 1996) or ancient DNA (Mathieson et al. 2015) where available. However, many traits do not leave a fossil record, and ancient DNA is not always available. An alternative approach is to examine the genomes of contemporary individuals, perhaps in combination with the estimated effects of contemporary variants on target traits (Berg and Coop 2014; Robinson et al. 2015; Field et al. 2016; Racimo et al. 2018; Uricchio et al. 2019), examining either allele-frequency differences among groups or traces of selection in patterns of within-population genetic diversity. If allelefrequency changes can be estimated from contemporary genetic data, then those allele-frequency changes, in combination with information about associations between alleles and traits of interest, can be used to study selection on traits. One reason estimated ARGs may be useful in population-genetic analysis of such subtle signals is that, to the extent the estimated ARG is correct, it naturally integrates information from flanking genomic regions in a way that reflects the history of recombination near the locus (Link et al. 2023). Relatedly, others have noted that even if estimated tree sequences are incorrect, the fact that they integrate information from nearby segments can prove useful in downstream inference (Whitehouse et al. 2024).

Edge and Coop (2019) proposed a set of methods to estimate the historical time course of a predicted population-mean trait level using local coalescent trees embedded in an ARG. The trait prediction is known as a polygenic score (PGS) or polygenic index. Although some of their methods are applicable to any trait prediction formed from genetic data, Edge and Coop focused on a population-mean PGS expressed as a weighted sum of population allele frequencies at unlinked loci, $Z(t) = 2 \sum_{i=1}^{k} \beta_i p_i(t)$, where the weight β_i is the additive effect size on a trait of interest of an allele at locus i, and $p_i(t)$ is the frequency of the effect allele at locus i at time t. The estimators operate by estimating allele-frequency changes at the loci contributing to the PGS. Under neutrality, the proportion of lineages at time t in a coalescent tree subtending derived alleles in the contemporary sample is an unbiased estimator of the derived-allele frequency at time t. Under selection, the ancestors of the sample are a biased sample from the ancestral population, but ideas from phylodynamics can be borrowed to form noisy estimates of the number of carriers of each allelic type in a specified time period.

To gain a deeper understanding of the performance of different ARG estimation approaches in estimating population-mean PGS histories, we applied the Edge and Coop framework to estimate local trees from 7 methods for ARG estimation. For comparison with the work of Edge and Coop (2019), we evaluated RENT+ (Mirzaei and Wu 2017), which they used in their original paper. We also evaluated the performance of ARGweaver (Rasmussen et al. 2014), which predates RENT+ and provides more accurate estimates on smaller samples, as well as Relate (Speidel et al. 2019), tsinfer+tsdate (Kelleher et al. 2019; Wohns et al. 2022), ARG-Needle and ASMC-clust (Zhang et al. 2023), and SINGER (Deng et al. 2024), all of which scale to larger samples. We also consider the effect of sample size on the resulting estimates.

Methods

We simulated derived-allele frequency trajectories at all unlinked loci that affect a trait. The true population-mean PGS trajectory was computed as the weighted sum $2\sum_{i=1}^{k} \beta_i p_i(t)$ of these allelefrequency trajectories, where the weights β_i are additive effect sizes, and $p_i(t)$ is the frequency of the effect allele at locus i at time t. Taking these simulated allele-frequency trajectories as input, we generated corresponding coalescent trees and haplotypes using mssel (Berg and Coop 2015). The haplotypes were fed to all of the tree estimation software packages considered here, and the estimated trees were saved. Upon obtaining both the true and estimated trees, we applied the allele-frequency estimators from Edge and Coop (2019) to them, and the estimated allele-frequency trajectories were used to compute estimated allele-frequency gestimated population-mean PGS trajectories using the additive effect sizes. The estimated population-mean PGS trajectory was then compared with the true trajectory using various metrics (Fig. 1).

Simulations

We simulated the population-mean PGS trajectories of traits additively determined by 100 unlinked loci under two scenarios: (i) neutral evolution and (ii) trait-increasing directional selection occurring from 0.04 to 0.02 coalescent units ago, with neutrality at other time points. Assuming an effective population size of 10,000 and a generation time of 30 years, the period of 0.02-0.04 coalescent units in the past corresponds to 12,000-24,000 years in the past, or the most recent part of the Upper Paleolithic. Because the bulk of ancient DNA evidence is from samples more recent than this (Mathieson et al. 2015; Speidel et al. 2021; Stern et al. 2021), indirect methods to detect selection in humans are perhaps especially of interest in this epoch, given that allelefrequency changes cannot as easily be examined directly (Le et al. 2022; Mathieson and Terhorst 2022). The PGS is a weighted sum of the allele frequencies, with weights equal to the additive effect sizes. There are no environmental effects or interactions-in our simulations, there is no distinction between the PGS and the trait. Further, we treat all the effect sizes as known.

For each trait, we simulated the derived-allele frequency histories of 100 trait-associated loci. For each locus, an effect size for the derived allele was drawn from a normal distribution $\mathcal{N}(0, h^2 \sigma^2 w/n)$ in which the heritability h^2 and the contemporary variance σ^2 of the trait are set to 1, w is a modified version of Watterson's constant calculated as $\sum_{i=\lceil 2Nc\rceil}^{\lfloor (1-c)2N \rfloor} 1/i$ (N is the effective population size; c is the minimum minor allele frequency being drawn, we set c = 0.01), and *n* is the number of loci affecting the trait, set to 100. We consider a hypothetical trait with a heritability of 1 and for which "true" additive effect sizes are known. We do this because we are interested in considering estimation accuracy at the level of a population-mean polygenic score rather than the trait to which the polygenic score corresponds. Polygenic scores will differ from trait values due to biases and errors in genomewide association study (GWAS) estimation and SNP heritabilities less than one, among other factors. Further, back in time, systematic changes in the environment and difficulties in "porting" polygenic scores trained in modern samples onto ancient samples will affect accuracy. Thus, we limit our focus to accuracy in population-mean polygenic score estimation. However, we emphasize that the model for selection we use assumes that selection occurs on the polygenic score itself. Given that most polygenic scores are noisy predictors of the traits with which they are associated, the selection gradients we simulate might be thought of as corresponding to larger selection gradients on actual trait values.

Next, we simulated allele-frequency trajectories for each locus. For both neutral and selected traits, the majority of the history of each locus was simulated backward in time—for neutral traits, simulation was entirely backward, and for selected traits, simulation was backward-in-time prior to the period in which selection occurred. The simulation is forward-in-time during the selection



Fig. 1. Method overview. a) Simulate derived-allele frequency trajectories at unlinked loci. b) Generate true trees and haplotypes from those trajectories. c) Obtain estimated trees from tree estimation software. d) Apply allele-frequency estimators to true and estimated trees, then compare the true and estimated population-mean PGS trajectories.

interval and afterward. We set the probability of obtaining a derived-allele frequency k/2N at the more recent end of the period simulated backward in time to be inversely proportional to k, where $k \in \{1, ..., 2N - 1\}$. This sampled allele frequency serves as the present derived-allele frequency under neutrality and as the initial allele frequency at the onset of selection in the recent-selection scenario.

Given a derived-allele frequency $p_i(t)$, the allele frequency at 1/(2N) coalescent time units in the past was drawn from $\mathcal{N}(p_i(t)(1-1/2N), p_i(t)(1-p_i(t))/2N)$ (Przeworski et al. 2005; Berg and Coop 2015; Lee and Coop 2017). During the selection period, we simulated the derived-allele frequency forward in time in steps of 1/(2N) coalescent units. The frequency $p_i(t + 1)$ was drawn from $\mathcal{N}(p_i(t) + sp_i(t)(1 - p_i(t)), p_i(t)(1 - p_i(t))/2N)$ conditional on the frequency $p_i(t)$, where s is the selection coefficient on the derived allele at time t. The value of s is $\alpha\beta$, where α is the selection gradient on the trait at time t and β is the effect size of the derived allele. (We set the value of α according to the expected trait variance at the onset of selection, which was 1.) This procedure is an approximation of allele-frequency dynamics under polygenic selection, but it is one that captures the overall patterns we seek to study (see supplementary text and Supplementary Fig. S1 for a comparison with truncation selection simulated forward in time).

After the allele-frequency trajectory was simulated, the polygenic-score trajectory calculated based on the allele-frequency trajectory was retained if the difference in population-mean PGS between the onset and end of the selection was within 5% of the expected change $2N\delta tS$, where δt is the duration of selection in coalescent units and S is the selection differential on the PGS. We

imposed this 5% cutoff in order to ensure that the degree of trait change during the period of selection was similar across iterations, and to ease comparisons with the results of Edge and Coop (2019), who used the same cutoff. It leads to selection of the 25–30% of simulated trait trajectories closest to the expectation in the parameter regime we used.

We used an unpublished modified version of ms (Hudson 2002) called mssel (Berg and Coop 2015) to produce simulated local coalescent trees and haplotypes across a region flanking the causal variant. In mssel, for most simulations, the sample size was set to 2,000 and the number of derived chromosomes was drawn as a binomial random variable with a size of 2,000 and success frequency equal to the contemporary allele frequency. We selected an effective population size N of 10,000, and a haplotype length of 200,000 base pairs (with the effect locus at position 100,000). The per-base-pair mutation rate was set as 2e-8, and the per-base-pair recombination rate was set as 2.5e-8. These values were transformed to population-scaled mssel inputs of -r 199.5 and -t 159.68. To explore additional scenarios, we also independently simulated larger sample sizes (5,000), haplotype of 500,000 base pairs, realistic human demography, genotyping error, and phasing error.

Software specifications

The coalescent trees produced by mssel for the selected sites are the true trees, and the haplotypes corresponding to the true trees were used as input to ARG estimation software to generate estimated trees. A brief description of each piece of estimation software can be found in supplementary text. The ARG estimation programs vary with respect to the total sample size that can be used. We thus down-sampled the 2,000 simulated haplotypes generated in each simulation as necessary. We used 20 randomly drawn haplotypes as input to ARGweaver, and 20 and 200 haplotypes for RENT+ and SINGER. For comparison with these methods, we also used the same subsamples of 20 and 200 haplotypes as input to Relate and tsinfer+tsdate. We did not use subsamples for ARG-Needle or ASMC-clust, which require sample sizes of at least 300 haplotypes.

ARGweaver

ARGweaver site input files were generated by a custom Python function. We ran the arg-sample program with the same values in the simulation (-N 10000, -r 2.5e-8, -m 2e-8) with the SMC' model (-smcprime). The total number of sample iterations (-n) was set to 700 with the first 200 iterations as burn-in, and 50 estimated trees were extracted for one locus, one from every 10th iteration (-sample-step 10).

RENT+

Haplotypes generated by mssel were converted to RENT+ format. We ran RENT+ using -t to estimate branch lengths for local trees and -1 to specify the proportional positions on the chromosomes.

Relate

We wrote a script to generate Relate haps, map, and sample files from simulated haplotypes. Relate was run with -mode All, -m 2e-8 (mutation rate), and -N 20,000 (haploid effective population size). We used the add-on module RelateExtract to extract the tree corresponding to the selected site.

tsinfer+tsdate

The tsinfer input file was generated by calling the add_site function. Coalescent trees were estimated with tsinfer using default settings. The age of nodes in the tree was estimated using tsdate with the same parameter values as in the simulations (Ne=10,000, mutation_rate=2e-8). We converted the tree at the selected site to Newick format with the as_newick method in tskit.

ARG-Needle and ASMC-clust

The haps, map, and sample input file of ARG-Needle were generated by custom R functions. We set the -normalize parameter as the default value and used the constant 20K-sized (haploid population size) demography for -normalize demography. ARG-Needle requires a "decoding file" containing information about the demographic model, time discretization, and allelefrequency information. We created our own decoding file with the prepare decoding function in Python package asmcderived-allele frequencies were calculated on the basis of the simulated haplotypes; the discretized time intervals were set as 14 intervals of 30 generations each starting in the present (0, 30, 60..., 420), followed by 14 more ancient intervals of 100 generations (520, 620, 720,..., 1920); the demography file was built with a constant population size ($N_e = 20,000$). The output tree was converted to Newick format with the arg to newick function from the arg needle lib package. Both choices for parameter --mode ("array" and "sequence") were tested. To run ASMC-clust, we set the parameter -asmc clust to 1. Other parameters were kept at the default values.

SINGER

The VCF file was generated from the simulated data by a custom R function. The mutation rate (-m) was set to 2e-8 and the ratio (-ratio) was set to 1.25. The population size (-Ne) was set to 10,000. SINGER was run for 7,000 iterations, with the first 2,000 iterations as burn-in and thinning every 10 iterations (-thin 10). We took 50 samples from the 5,000 posterior trees with interval size 100 (-step 10). The estimated trees were stored as tree sequences. The marginal tree at the selected site was converted to Newick format with the as_newick method in tskit.

Branch length units

The ARG estimation tools we considered use different units for branch length. We standardized branch lengths to units of 2N generations for all trees. True trees from mssel are produced in units of 4N generations, so we multiplied the branch lengths by 2. In the manuscript describing RENT+, there is some ambiguity about the units, which are described as "standard coalescent units" (typically units of 2N generations), but in some calculations appear as if they are in units of 4N generations (Supplementary Tables S1 and S2). We applied the estimators on both the original branch lengths (assuming units of 2N generations) and a rescaled branch length (assuming the reported values are in units of 4N generations). We show the rescaled version in comparison with the other methods—although Edge and Coop (2019) treated the reported branch lengths as if they were in units of 2N generations, we believe that the makers of RENT+ intended the branch lengths to be interpreted as in units of 4N generations. However, assuming 2N-generation units leads to slightly better performance on average (Supplementary Figs. S2 and S3, Table S1). The ARGweaver, tsinfer+tsdate, ARG-Needle, ASMC-clust, and SINGER trees are reported in units of generations, and we divided their original branch length by 2N. Relate trees use years as the unit and assume 28 years per generation when exported in Newick format, so we divided the original branch length by 28 × 2N. Average times to most recent common ancestor (tMRCA) from the original branch length and the scaled branch length for each method are listed in Supplementary Tables S1 and S2.

Applying population-mean PGS trajectory estimators to the estimated trees

We applied the three estimators proposed in Edge and Coop (2019): the "proportion-of-lineages" estimator, the "waiting-time" estimator, and the "lineages-remaining" estimator (further description of the estimators in supplementary text). All three estimators were applied to the estimated trees to obtain the estimated allele-frequency trajectories for each locus. For ARGweaver and SINGER trees, since there is more than one estimated tree for each locus, we applied the estimators to each sampled tree and took the average result from 50 trees as the estimated frequency trajectory for that locus. As discussed in Edge and Coop (2019), the proportion-of-lineages estimator is expected to perform well under neutral evolution but to be biased during periods of directional selection. The other two estimators are expected to be approximately unbiased given the true trees but to be much more variable. The waiting-time estimator, as written in Edge and Coop (2019) is not applicable in cases in which there are polytomies in a coalescent tree. We devised a scheme to address this limitation (see supplementary text).

As a general note about the estimators of Edge and Coop (2019), we emphasize that the relationship between changes in the population-mean PGS and trait changes is not straightforward. Even if PGS histories are estimated perfectly on an accurate PGS, phenotypes may be influenced by environmental changes and gene-by-gene or gene-by-environment interactions. Additionally, the linkage disequilibrium (LD) between markers and causal variants will change over time, causing the relationship between markers and the phenotype to change with it (Martin *et al.* 2017). Finally, some variants affecting the trait in the past will have been lost in the present (Carlson *et al.* 2022) or may no longer be detectable in GWAS. Nonetheless, Edge and Coop's methods can be used to identify and roughly date patterns of allele-frequency change that are not consistent with neutral evolution.

As noted by Edge and Coop (2019), like many methods for detecting selection from GWAS information, their methods are vulnerable to biases from population stratification or assortative mating, which can create spurious signals of directional selection. One potential countermeasure is to use effect size estimates from family-based studies, which are less sensitive to such biases, though not completely immune (Veller and Coop 2024; Veller et al. 2024).

Benchmarking metrics

We evaluated the performance of the ARG-estimation software as inputs to various approaches to PGS-history estimation, comparing them in terms of bias, mean squared error (MSE), 95% confidence interval coverage, and type I error rates and power (at the 0.05 level) of the T_X statistic for testing directional selection (Edge and Coop 2019). The T_X statistic is computed as a sum of squared, standardized changes in an estimated population-mean PGS between prespecified time points. That is, $T_X = \sum_{i=0}^{w} X_i^2$, with

$$X_{j} = \frac{Z(t_{j}) - Z(t_{j-1})}{\sqrt{2V_{A}(t_{j} - t_{j-1})}}$$

where $Z(t_j)$ is the population-mean PGS at time t_j , and V_A is the additive genetic variance of the PGS $(2\sum_{i=1}^k \beta_i^2 p_i(1-p_i))$. Under neutral evolution, X_j approximately follows a standard normal distribution, and the allele-frequency changes are independent in distinct time intervals. Thus T_X , the sum of the X_j^2 , should approximately follow the $\chi^2(w)$ distribution. However, under directional selection, the population-mean PGS changes more quickly than predicted under neutrality, leading to T_X values that are large compared with the expected $\chi^2(w)$ distribution. T_X is sensitive to directional selection. Although Appendix A of Edge and Coop (2019) is suggestive of a modification of T_X that might be sensitive to stabilizing selection, no such statistic has yet been proposed or studied in this framework.

The efficacy of Edge and Coop's estimators depends on the accuracy of branch lengths in the estimated trees. To explore the differences in software performance, we extracted the pairwise coalescent times from true trees and estimated trees. Then, we compared the point estimate of coalescent times between the two, and the overall distribution of pairwise coalescent times from estimated trees against the expected exponential distribution, similar to the benchmarking recently performed by Brandt *et al.* (2022). To obtain the samples of pairwise coalescent times, we gathered trees containing the causal allele from 10 randomly selected simulated traits, resulting in a total of 1,000 trees. Specifically, for Relate, tsinfer, ARG-Needle, and ASMC-clust, we collected pairwise coalescence times from their 2,000-sample trees. For RENT+ and SINGER, we used 200-sample subtrees, and **Table 1.** Average runtime for estimating ARGs for one 200 kb segment of 20, 200, 2,000 and 5,000 samples from the ARG-estimation tools.

Sample size	20	200	2,000	5,000
RENT+	3.9	38	_	_
ARGweaver	2,723	_	_	_
Relate	0.4	2.5	392.9	1,863.5
tsinfer	0.9	2.2	6.7	37.1
SINGER	259	917	_	_
ARG-Needle	_	_	165	422.6
ASMC-clust	_	_	1,594	-

Runtimes are measured in seconds.

for ARGweaver, we used 20-sample subtrees. In the case of ARGweaver and SINGER, we averaged pairwise coalescence time across 50 sampled trees per locus. The averaged times correspond to the same pairs of tips in each sampled tree. Due to memory constraints, we could not plot all pairwise coalescence times from one thousand distinct 200-sample and 2,000-sample trees. Instead, we randomly sampled 190,000 pairwise coalescence times, which is 1,000 times the total number of pairs of tips on a 20-tip tree (i.e. 1,000 times $\binom{20}{2}$).

We also examined the impact of a selection event on the distribution of pairwise coalescence times by setting the selection coefficient to ~0.004, causing the derived-allele frequency to increase from approximately 0.3 to 0.7 during the selection period. Then, 100 estimated trees with 300 samples were generated on the basis of this allele-frequency trajectory. We compared the MSE of the estimates from true trees with different sample sizes (20, 200, and 2,000 haplotypes) to assess the effect of the increased sample size.

To explore the relationship between measures of tree-topology accuracy and performance in estimation of population-mean PGS history, we examined the Robinson–Foulds (Robinson and Foulds 1981) and Kendall–Colijn distance (Kendall and Colijn 2016) between true trees and estimated trees and the proportion of estimated trees with monophyletic derived tips.

Results

Runtime comparison

The ARG-estimation tools we consider vary substantially in their runtime. One important difference is between tools that incorporate MCMC sampling (ARGweaver and SINGER) and those that do not. Both ARGweaver and SINGER sample trees from the posterior distribution, and the reported runtime represent 700 MCMC samples (200 burn-in) from ARGweaver and 7,000 MCMC samples (2,000 burn-in) from SINGER. SINGER is, as expected, substantially faster than ARGweaver. Additionally, tsinfer has a clear advantage in speed with larger sample sizes. Average runtimes for every tool are listed in Table 1.

Bias, mean squared error, and confidence interval coverage under neutrality

When traits are simulated under neutral evolution and the true trees are known, the proportion-of-lineages procedure can be viewed as a maximum-likelihood estimator for the allelefrequency history (Edge and Coop 2019). The waiting-time and lineages-remaining estimators have previously been observed to be more variable than the proportion-of-lineages estimator with either the true trees or RENT+ trees (Edge and Coop 2019).



Fig. 2. The bias (a–c), MSE (d–f), and confidence interval coverage (g–i) of the proportion-of-lineages (left column), waiting-time (middle column), and lineages-remaining estimators (right column), with the true trees and estimated trees from each ARG-estimation method as input. For the estimates computed from true coalescent trees (black lines), 1,000 simulations were performed with a sample of 2,000 chromosomes. In each simulation, the PGS was formed from 100 loci and evolved neutrally. Relate, tsinfer, and ASMC-clust were used to reconstruct trees with 2,000 chromosomes. SINGER was run with 200 chromosomes, and ARGweaver was run with 20 chromosomes. Lines represent means from 100 simulations. Times are displayed assuming diploids with N_e = 10,000 and a generation time of 30 years, i.e. one coalescent unit corresponds to 600,000 years.

With respect to bias, we find that, in line with previous results, none of the estimators display much bias when estimated with the true trees or RENT+ under neutrality. In addition, almost all newly tested tools/algorithms (Relate, tsinfer+tsdate, ARGweaver, ASMC-clust, and SINGER) follow similar patterns. Differences in bias resulting from trees estimated with these different pieces of software, especially during the recent past, and especially with the proportion-of-lineages estimator, are relatively minor (Fig. 2a-c, Supplementary Figs. S4-S6). Results for ARG-Needle are reported in the supplement and show more bias than other estimated trees (Supplementary Fig. S6). The likely reason for this is that ARG-Needle's algorithm does not guarantee that mutations map to particular branches in the ARG. Instead, carriers of the derived allele may form polyphyletic groups on the local marginal tree. The frequency of nonmonophyly among tips carrying the derived allele is much higher in

ARG-Needle trees than in trees estimated by any other method (Supplementary Table S3). Although such polyphyly may be acceptable for many purposes, it can cause major problems for the estimators of Edge and Coop (2019), particularly in the ancient past. However, if ARG-Needle trees are manually modified to force monophyly among derived-allele-carrying tips, then the bias of ARG-Needle is similar to other methods (Supplementary Fig. S7). We also found that when the "-mode" parameter in ARG-Needle is set to "sequence," the results improve over the default value of "array" (Supplementary Fig. S7), which is sensible given that the simulated data we provide to ARG-Needle includes all variants flanking the focal site.

The advantage of the proportion-of-lineages estimator over the waiting-time and lineages-remaining estimators under neutrality is more apparent when examining estimated MSE. In the three MSE plots (Fig. 2d–f), MSE tends to increase from the recent past

into the distant past for all estimators and all software packages. This is unsurprising. Under neutrality, the variance of the proportion-of-lineages estimator is inversely related to the number of lineages ancestral to the sample (Edge and Coop 2019). Thus, the fact that most lineages coalesce in the recent past implies that the estimator is more variable in the distant past. The proportion-of-lineages estimator (Fig. 2d) exhibits a much lower MSE than the other two estimators. The true trees, again unsurprisingly, provide the best MSE. However, among the estimated ARGs, the ones produced by software that scales to large samples do not necessarily show markedly better performance. Across much of the range examined, SINGER, which uses samples of only 200 haplotypes, has the lowest MSE, and in the more distant past, ARGweaver, which uses samples of only 20 haplotypes, is comparable with ASMC-clust, which uses samples of 2,000 haplotypes. In the MSE plots of the waiting-time and lineages-of-remaining estimators, the estimates derived from ARGweaver and SINGER trees consistently show relatively low MSE values (Fig. 2e, f).

We assessed credible interval coverage with SINGER and ARGweaver and confidence interval coverage with the true trees and all other methods (Fig. 2g-i). The true trees produce acceptable coverage with all estimators. Credible intervals from SINGER-estimated trees also show consistently high coverage in most cases, though they drop below 95% coverage in the recent past with the proportion-of-lineages estimator. Other estimated trees, however, all produced somewhat lower coverage, with a tendency for coverage to decline into the more distant past.

ARGweaver and SINGER sometimes outperform the true trees in terms of MSE or credible interval coverage when the lineagesremaining and waiting-time estimators are used, a counterintuitive result. However, individual ARGweaver and SINGER trees do not outperform the true trees (Supplementary Figs. S8 and S9); it is only when results are averaged across the many trees produced by these methods that performance exceeds the true trees. The lineages-remaining and waiting-time estimators are noisy estimators that are sensitive to the timing of individual coalescent events, and this result suggests that their variability can be reduced by averaging results across many trees compatible with the data.

The bias, MSE, and coverage under neutrality with larger flanking regions, a nonconstant demography, simulated genotyping error, and simulated phasing error can be found in Supplementary Figs. S10–S13. (We did not run ARGweaver or ASMC-clust in these conditions because of their longer runtimes.) The results were similar to those in Fig. 2 overall, though Relate and tsinfer +tsdate performance under the waiting-time and lineagesremaining estimators appeared to decline somewhat with genotyping error and a more realistic human demography. Performance on the proportion-of-lineages estimator, which is preferred under neutrality, remained about the same.

Bias, mean squared error, and confidence interval coverage under recent directional selection

Under selection, the proportion-of-lineages estimator is known to be biased, as the ancestors of the sample are not representative of the ancestral population from which they are drawn (Edge and Coop 2019). The waiting-time and lineages-remaining estimators were developed to avoid this bias. As pointed out previously (Edge and Coop 2019), because different estimators perform well under neutrality and directional selection, one reasonable procedure is to test for selection using T_X and then choose an estimator on the basis of the result. When there is a burst of directional selection in the recent past, the performance differences among ARG-estimation software packages are more pronounced. In contrast to the neutral scenario, and as expected, the proportion-of-lineages estimator is more strongly biased, has a higher MSE, and has lower confidence interval coverage than the other estimators (Fig. 3a–i, Supplementary Figs. S14–S16), particularly during the period of selection and more anciently, as expected. Looking backward in time, the bias of the proportion-of-lineages estimator between the present and the end of the period in which selection occurred is low. This basic pattern also appears in the MSE and confidence interval coverage plots, with good performance between the present and the end of directional selection, declining into the past during the period in which selection occurred, and then recovering during the neutral period that preceded selection.

The waiting-time and the lineages-remaining estimators show substantially lower bias (Fig. 3b and c). The true trees generally show acceptable performance throughout the time period examined, with only slight bias during the period of selection and fairly uniform performance across time on other desiderata. However, all estimated trees produce noticeably worse performance on all criteria (Fig. 3b–i). Overall, among the estimated trees, the SINGER trees produce the lowest bias and MSE and interval coverage closest to the nominal level. ARGweaver ranks second across most of the investigated time. ASMC-clust and Relate are competitive in the recent and distant past respectively.

As in the neutral case, we also compared performance between individual trees and averages across many posterior ARGweaver and SINGER trees (Supplementary Figs. S17 and S18). We also compared the performance of original and modified ARG-Needle trees (Supplementary Fig. S19). We examined the performance of SINGER, Relate, and tsinfer+tsdate for scenarios with larger flanking regions, simulated genotyping error, and simulated phasing error (Supplementary Figs. S20–S23). Analogously to the neutral case, these changes led to broadly similar results to those in Fig. 3, though the performance of Relate and tsinfer+tsdate declined somewhat under genotyping error and under a realistic human demography.

Power of T_X

Next, we examined the type I error rate and power of tests of neutrality using T_X , a test statistic sensitive to changes in a population-mean PGS that are larger than expected under neutrality. T_X can be understood as a version of the Q_X statistic (Berg and Coop 2014) applied to a population-mean PGS from one population through time, rather than multiple populations sampled at the present (Edge and Coop 2019). To use T_X , one picks a set of time points to calculate the statistic $X_j = \frac{Z(t_j)-Z(t_{j-1})}{\sqrt{2V_A(t_j-t_{j-1})}}$, where $Z(t_j)$ is the estimated population-mean PGS at time t_j and V_A is the additive genetic variance of the PGS. With the true allele-frequency trajectories, under neutrality, the sum across time points $T_X = \sum_{j=1}^{w} X_j^2$ computed from w distinct intervals is approximately $\chi^2(w)$ distributed (Edge and Coop 2019).

In line with previous results, T_X has an acceptable type I error rate when compared against the χ^2 distribution only when the proportion-of-lineages estimator is used to form the statistic (Table 2, Supplementary Table S4). With the proportion-of-lineages estimator, the observed type I error rates for the true trees, tsinfer+tsdate trees, ASMC-clust trees, and SINGER trees do not differ significantly from the nominal rate (RENT+ and ARG-Needle trees in Supplementary Tables S5 and S6).



Fig. 3. The bias (a–c), MSE (d–f), and confidence interval coverage (g–i) of the proportion-of-lineages (left column), waiting-time (middle column), and lineages-remaining estimators (right column), with the true trees and estimated trees from each ARG-estimation method as input. The PGS is influenced by a period of directional selection that occurred from 0.04 to 0.02 coalescent units ago as represented by the vertical dotted lines. The parameters used to run all software are otherwise identical to those applied under neutrality (Fig. 2).

Input		Proportion-of-lineages		
		χ^2 distribution	Permutation distribution	
Neutral	True trees	0.069	0.054	
	Relate	0.1*	0.07	
	tsinfer	0.07	0.04	
	ASMC-clust	0.04	0.02	
	SINGER	0.09	0.05	
Selection	True trees	0.982	0.981	
	Relate	0.14	0.21	
	tsinfer	0.19	0.07	
	ASMC-clust	0.69	0.57	
	SINGER	0.9	0.85	

For the type I error simulations, asterisks indicate whether the observed type I error rate differs significantly from the nominal rate of 0.05: * P < 0.05.

For Relate, the type I error rate is slightly higher than nominal. However, with all estimators and tree estimation software, calibrated type I error rates can be recovered by using a permutation distribution instead of the theoretical χ^2 distribution, as expected.

We also tested T_X 's power in simulations that included a period of directional selection between 0.02 and 0.04 coalescent units in the past. In line with the previous results, T_X calculated from the proportion-of-lineages-estimated allele frequencies is much more powerful than when the other two estimators of allele frequency are used. Considering proportion-of-lineages T_X , power was approximately 98% with the true trees. Except for SINGER, most estimated trees lead to a substantial loss in power. Comparing against the permutation distribution, SINGER trees led to an observed power of 85%, whereas ASMC-clust produced T_X statistics with a power of 57%, and Relate and tsinfer+tsdate trees



Fig. 4. The comparison of pairwise coalescence times from true trees and estimated trees under directional selection. Each dot represents the true and estimated (log) coalescence time between a pair of samples. The diagonal line shows x = y, i.e. true times equal to estimated times. The values in the top-left corner show Spearman correlation coefficients.

produced power estimates of 21 and 7%, respectively. Power obtained with RENT+ trees with different branch length units can be found in Supplementary Table S5.

Comparison between true and simulated pairwise coalescence time

A major contributor to differences in performance among the tree-estimation procedures is the accuracy of the estimated coalescence times. We compared the pairwise coalescence times from estimated trees with their true values. As measured by the Spearman correlation between true and estimated pairwise coalescence times, SINGER outperforms other software packages, with ARGweaver and ASMC-clust close behind. It is also possible to see finer-grained patterns in the log-scaled coalescence times displayed in Fig. 4. Overall, all three of SINGER, ARGweaver, and ASMC-clust show a tendency to overestimate coalescence times. But for short coalescence times, ASMC-clust tends toward underestimation (Fig. 4). Both Relate and tsinfer+tsdate estimates appear biased for moderate-length coalescence times (between 0.1 and 1); Relate tends to underestimate these values, whereas tsinfer+tsdate tends to overestimate them (Fig. 4). Additionally, Relate produces more variable estimates of short coalescence times than tsinfer+tsdate, as previously observed by Fan et al. (2023). Finally, the time discretizations used by ARGweaver and ASMC-clust are visible as horizontal bands. These patterns are qualitatively similar when examined in simulations performed under neutrality (Supplementary Fig. S24, see also Brandt et al. 2022) and in general distributions of pairwise coalescence times (Supplementary Figs. S25 and S26).

Supplementary Table S7 shows the correlation between the true and estimated time to most recent common ancestor (tMRCA) for each software package and sample size under neutrality. With the exception of RENT+, the correlation coefficients stay stable across varying sample sizes for Relate, tsinfer +tsdate, and SINGER.

To examine the effect of selection on coalescence time estimation, we checked the distribution of pairwise coalescence time from 100 trees with 300 tips, each subject to a selection event that increases the allele frequency from 0.3 to 0.7 between 0.04 and 0.02 coalescent units in the past. This selection event results in a high density of coalescences in the recent past within the derived-allele subtree, leading to a higher frequency of short coalescent times compared with the neutral simulations. This creates a distinct dip in the histogram of pairwise coalescence times (Fig. 5). This pattern can also be observed, to varying degrees, in the distributions of pairwise coalescence times from estimated trees (Fig. 5). Nevertheless, the dip in the histogram is wider and deeper when examining the true trees than when examining times from estimated trees. The density of recent coalescence times is also too low in estimated trees.

Topological aspects of ARG estimation accuracy

In addition to the correlation between true and estimated pairwise coalescence times, we considered other aspects of ARG estimation accuracy. Specifically, with respect to topology, we considered the rate of polyphyly among derived tips (and related measures), the Robinson–Foulds distance, and the Kendall–Colijn distance (with $\lambda = 0$, i.e. excluding branch length information). Results are shown in Supplementary Tables S3 and S8.

Each of these measures appears to track performance with respect to the estimators we consider here to some degree, but there are observations in which software that produces more accurate trees via each topological metric produces less accurate population-mean PGS histories. For example, with 200 tips, RENT + produces trees with lower Robinson–Foulds distances than SINGER, despite SINGER's excellent performance on our benchmarks. With respect to the Kendall–Colijn distance, with 200 tips, tsinfer trees perform better than both Relate and RENT+, despite the latter's generally better performance on our benchmarks.

Impact of sample size on estimation results

Another potential contributor to differences in performance among estimates derived from distinct ARG estimation



Fig. 5. Distribution of pairwise coalescence times from true and estimated trees at a locus undergoing strong selection. Histograms show times from 100 trees with 300 tips each. Each tree (and flanking regions) was simulated assuming strong selection that increases the minor allele frequency from 0.3 to 0.7 between 0.04 and 0.02 coalescent units in the past.

procedures is sample size—ARG estimators differ in the sample sizes they can accommodate, and this variation is reflected in our simulations. Because large samples of haplotypes coalesce quickly in the recent past, we expect increased sample size to benefit allele-frequency trajectory estimation primarily in the recent past. We compared estimated PGS histories formed from the true trees obtained from samples of different sizes. The MSE of the proportion-of-lineages estimates decreases substantially when the sample size increases from 20 to 200. Notably, the MSE is not much larger for samples of size 200 than for samples of 2,000 lineages under neutrality, when this estimator is expected to perform well (Fig. 6a). Even in the recent past, the difference is slight. This pattern extends to the other two estimators when applied under selection (Fig. 6e and f). But under neutrality, the waiting-time and lineages-remaining estimators perform approximately equally regardless of sample size (Fig. 6b and c). The pattern observed for true trees also applies to estimated trees (Supplementary Figs. S4-S6, S14-S16).

Although we see little improvement in MSE when using true trees with samples of 2,000 relative to samples of 200 haplotypes, it is possible that large samples aid in marginal tree estimation, such that large samples produce better results with methods that can handle their size, not because the additional tips are helpful per se, but rather because the trees they produce are more nearly accurate. We do not find clear evidence of this possibility when increasing to 2,000 or 5,000 samples (Supplementary Figs. S27 and S28), but tsinfer+tsdate and ARG-Needle can accommodate samples much larger than this, which we do not explore.

Empirical data analysis

Edge and Coop (2019) analyzed data from the GBR (British) subsample of the 1000 Genomes project with respect to polygenic predictions of height formed from effect sizes from either GIANT (Wood *et al.* 2014) or the UK Biobank (Neale Lab 2017). They found that GIANT effect sizes produced an impression of an increase in height over the last 60–90 ky in ancestors of the GBR individuals, but UK Biobank effect sizes produced no such effect, consistent with recent (at the time) evidence that GIANT effect sizes were subject to biases from population stratification, small at the level of individual loci but large when combined into a polygenic score (Berg *et al.* 2019; Sohail *et al.* 2019). We repeated the analyses of Edge and Coop using trees estimated by Relate, tsinfer +tsdate, and SINGER (supplementary text and Supplementary Figs. S29–S31). The results are qualitatively similar to those of Edge and Coop (2019) in each case.

Discussion

We studied the performance of a set of methods for estimating the history of population-mean PGS using estimated ARGs, with a particular interest in the relative performance of different ARG estimation procedures. We used a broad range of methods appearing over the past decade—ARGweaver (Rasmussen et al.



Fig. 6. The MSE of the proportion-of-lineages (left column), waiting-time (middle column), and lineages-remaining estimators (right column) estimates from true trees with different sample sizes under neutrality (a-c) and directional selection (d-f).

2014), RENT+ (Mirzaei and Wu 2017), Relate (Speidel et al. 2019), tsinfer+tsdate (Kelleher et al. 2019; Wohns et al. 2022), ARG-Needle/ASMC-clust (Zhang et al. 2023), and SINGER (Deng et al. 2024)—that vary in their approaches, runtimes, and scalability. Previous work on these methods considered only estimated marginal trees from RENT+ (Edge and Coop 2019). We also benchmarked the estimated coalescence times emerging from these methods in general terms under neutrality and selection, providing an update on the work of Brandt et al. (2022), including several ARG estimation procedures that are newly released since their work.

Many of the basic patterns we observed were consistent across all ARG or marginal tree estimation procedures, as well as being consistent with previous work. For example, regardless of the tree estimation procedure used, the resulting PGS-history estimates are more accurate under neutral evolution than under selection, and PGS estimates designed to work either under neutrality ("proportion-of-lineages") or under selection ("waitingtime" or "lineages-remaining") had the expected pattern of relative performance.

At the same time, there were considerable differences in the accuracy of estimated population-mean PGS histories depending on the method used for marginal tree estimation. This is expected given the rapid development of ARG and marginal tree estimation methods in the past 10 years, which have resulted in increased scalability of several orders of magnitude (Brandt *et al.* 2022; Lewanski *et al.* 2024). However, bigger samples are not always better. In fact, the best performance overall from any set of estimated trees—whether measured in terms of mean squared error or confidence/credible interval coverage of the estimated PGS histories, or in terms of power in tests of natural selection—came from SINGER trees estimated with 200 haplotypes, rather than any of the tree estimation procedures that were capable of scaling to 2,000 haplotypes. There were several other cases in which smaller samples fit by either RENT+ or ARGweaver outperformed trees estimated from much larger samples.

The fact that trees estimated with much larger samples of haplotypes do not always outperform trees estimated with smaller samples may be counterintuitive. One part of the explanation is the tradeoff between accuracy and scalability in ARG estimation. ARG estimation is a difficult problem, and scalability can be achieved with simplifications that can reduce accuracy (Deng et al. 2021; Brandt et al. 2022, 2024; Deng et al. 2024). In line with this, ARGweaver and SINGER trees tend to feature more accurate branch lengths than methods that scale to larger samples (Brandt et al. 2022; Deng et al. 2024; Lewanski et al. 2024). The other important part of the explanation with respect to the estimators we explore here is that in large samples, coalescence initially happens very fast. Because the coalescence rate is proportional to $\binom{n}{2}$, with *n* the number of lineages that have not yet coalesced, large samples imply large amounts of coalescence in the very recent past. As a result, even very large samples will be represented by a small number of lineages in the recent past (Griffiths 1984; Tavaré 1984; Slatkin and Rannala 1997; Volz et al. 2009; Frost and Volz 2010; Maruvka et al. 2011; Chen and Chen 2013; Jewett and Rosenberg 2014), and therefore increasing the sample size beyond a few hundred haplotypes produces more precise allelefrequency estimation primarily in the very recent past. This is reflected in Fig. 6, which shows that, even with the true trees, increasing the sample size from 200 to 2,000 does not clearly reduce the MSE of the Edge and Coop's 2019 estimators in the scenarios in which they are each predicted to work well.

The ARG estimation methods we have examined have all been benchmarked previously. However, most of this benchmarking has been done with respect to general indicators of performance, such as the overall accuracy of local tree topologies (Rasmussen et al. 2014; Kelleher et al. 2019); a generalized Robinson–Foulds distance for ARGs (Zhang et al. 2023), the distribution of distances between topologically distinct local trees (Deng et al. 2021), and the distribution of pairwise coalescent times (Brandt et al. 2022; Deng et al. 2024). Although these general features are very important, it is not always straightforward to predict from them which methods will perform best when applied to a specific downstream task. For the estimators we explore here, at a given locus, all of them can be computed from the sets of coalescence times on the derived-allele and ancestral-allele subtrees. Thus, for these estimators, it is important that branch lengths are estimated accurately, as the branch lengths determine the coalescence times. Topology does not matter per se, but it does matter in practice, in that certain kinds of topological errors make it unlikely that the estimated coalescence times on each background will be nearly accurate. In line with this, we observe that SINGER excels in both our benchmarks and in preserving monophyly among tips carrying the derived allele, as expected in an infinite-sites mutational model.

For other tasks, other specific features may assume outsize importance. For example, in estimating the expected genetic relatedness matrix (eGRM), the estimation of long coalescence times matters a great deal, since many pairs of lineages, and thus many entries of the eGRM, are related by these long times, leading Relate to outperform tsinfer+tsdate for this task (Fan et al. 2022). In contrast, for some demographic inference problems, more recent coalescence times are more important, leading tsinfer+tsdate to outperform Relate (Fan et al. 2023). In addition to the specific features of the ARG that are important for a given downstream task, performance may vary according to the evolutionary scenario, e.g. the selection regime. In other words, there will not necessarily be an overall "best" method for ARG estimation-when choosing a method for ARG inference in empirical work, the specific strengths and weaknesses of particular methods may be more important than general considerations about overall accuracy or scalability.

This study adds to others suggesting that estimated ARGs and the marginal trees they encode are useful tools for the study of complex traits (Edge and Coop 2019; Chen and Chiang 2021; Speidel *et al.* 2021; Stern *et al.* 2021; Link *et al.* 2023; Zhang *et al.* 2023). As ARG estimation methods continue to improve in both accuracy and scalability, they will open new opportunities for mapping genetic variants contributing to phenotypes, revealing polygenic adaptation, and exploring the relationship between natural selection and population history.

Data availability

Code and the compiled version of mssel used in this article can be found at https://github.com/dandanpeng/ARG_Benchmarking. Supplemental material available at GENETICS online.

Acknowledgments

We thank members of the Edge, Mooney, and Pennell labs for helpful conversations and the anonymous peer reviewers for useful comments on the manuscript.

Funding

This work was supported by NIH grant R35GM137758 to MDE.

Conflicts of interest

The author(s) declare no conflicts of interest.

Literature cited

- Berg JJ, Coop G. 2014. A population genetic signal of polygenic adaptation. PLoS Genet. 10(8):e1004412. https://doi.org/10.1371/ journal.pgen.1004412
- Berg JJ, Coop G. 2015. A coalescent model for a sweep of a unique standing variant. Genetics. 201(2):707–725. https://doi.org/10. 1534/genetics.115.178962
- Berg JJ, Harpak A, Sinnott-Armstrong N, Joergensen AM, Mostafavi H, Field Y, Boyle EA, Zhang X, Racimo F, Pritchard JK, et al. 2019. Reduced signal for polygenic adaptation of height in UK Biobank. Elife. 8:e39725. https://doi.org/10.7554/eLife.39725
- Brandt DYC, Huber CD, Chiang CWK, Ortega-Del Vecchyo D. 2024. The promise of inferring the past using the ancestral recombination graph. Genome Biol Evol. 16(2):evae005. https://doi.org/10. 1093/gbe/evae005
- Brandt DYC, Wei X, Deng Y, Vaughn AH, Nielsen R. 2022. Evaluation of methods for estimating coalescence times using ancestral recombination graphs. Genetics. 221(1):iyac044. https://doi.org/ 10.1093/genetics/iyac044
- Carlson MO, Rice DP, Berg JJ, Steinrücken M. 2022. Polygenic score accuracy in ancient samples: quantifying the effects of allelic turnover. PLoS Genet. 18(5):e1010170. https://doi.org/10.1371/journal. pgen.1010170
- Chen H, Chen K. 2013. Asymptotic distributions of coalescence times and ancestral lineage numbers for populations with temporally varying size. Genetics. 194(3):721–736. https://doi.org/10.1534/ genetics.113.151522
- Chen M, Chiang CWK. 2021. Allele frequency differentiation at height-associated SNPs among continental human populations. Eur J Hum Genet. 29(10):1542–1548. https://doi.org/10.1038/ s41431-021-00938-2
- Christ R, Wang X, Aslett LJ, Steinsaltz D, Hall I. 2024. Clade distillation for genome-wide association studies. bioRxiv. pp. 2024–09. https://doi.org/10.1101/2024.09.30.615852, preprint: not peer reviewed.
- Deng Y, Nielsen R, Song YS. 2024. Robust and accurate Bayesian inference of genome-wide genealogies for large samples. bioRxiv. https://doi.org/10.1101/2024.03.16.585351, preprint: not peer reviewed.
- Deng Y, Song YS, Nielsen R. 2021. The distribution of waiting distances in ancestral recombination graphs. Theor Popul Biol. 141(11):34–43. https://doi.org/10.1016/j.tpb.2021.06.003
- Edge MD, Coop G. 2019. Reconstructing the history of polygenic scores using coalescent trees. Genetics. 211(1):235–262. https://doi.org/10.1534/genetics.118.301687
- Fan C, Cahoon JL, Dinh BL, Ortega-Del Vecchyo D, Huber C, Edge MD, Mancuso N, Chiang CWK. 2023. A likelihood-based framework for demographic inference from genealogical trees. bioRxiv. https://doi.org/10.1101/2023.10.10.561787, preprint: not peer reviewed.
- Fan C, Mancuso N, Chiang CWK. 2022. A genealogical estimate of genetic relationships. Am J Hum Genet. 109(5):812–824. https:// doi.org/10.1016/j.ajhg.2022.03.016
- Field Y, Boyle EA, Telis N, Gao Z, Gaulton KJ, Golan D, Yengo L, Rocheleau G, Froguel P, McCarthy MI, et al. 2016. Detection of human adaptation during the past 2000 years. Science. 354(6313): 760–764. https://doi.org/10.1126/science.aag0776
- Frost SDW, Volz EM. 2010. Viral phylodynamics and the search for an 'effective number of infections'. Philos Trans R Soc Lond B Biol Sci. 365(1548):1879–1890. https://doi.org/10.1098/rstb.2010.0060
- Griffiths RC. 1984. Asymptotic line-of-descent distributions. J Math Biol. 21(1):67–75. https://doi.org/10.1007/BF00275223

- Griffiths RC, Marjoram P. 1996. Ancestral inference from samples of DNA sequences with recombination. J Comput Biol. 3(4):479–502. https://doi.org/10.1089/cmb.1996.3.479
- Gunnarsson ÁF, Zhu J, Zhang BC, Tsangalidou Z, Allmont A, Palamara PF. 2024. A scalable approach for genome-wide inference of ancestral recombination graphs. bioRxiv. pp. 2024–08. https://doi.org/10. 1101/2024.08.31.610248, preprint: not peer reviewed.
- Harris K. 2019. From a database of genomes to a forest of evolutionary trees. Nat Genet. 51(9):1306–1307. https://doi.org/10.1038/ s41588-019-0492-x
- Harris K. 2023. Using enormous genealogies to map causal variants in space and time. Nat Genet. 55(5):730–731. https://doi.org/10. 1038/s41588-023-01389-9
- Hudson R. 1990. Gene genealogies and the coalescent process. Oxford Surveys Evol Biol. 7:1–44.
- Hudson RR. 2002. Generating samples under a Wright-Fisher neutral model of genetic variation. Bioinformatics. 18(2):337–338. https:// doi.org/10.1093/bioinformatics/18.2.337
- Jewett EM, Rosenberg NA. 2014. Theory and applications of a deterministic approximation to the coalescent model. Theor Popul Biol. 93:14–29. https://doi.org/10.1016/j.tpb.2013.12.007
- Kappelman J. 1996. The evolution of body mass and relative brain size in fossil hominids. J Hum Evol. 30(3):243–276. https://doi. org/10.1006/jhev.1996.0021
- Kelleher J, Wong Y, Wohns AW, Fadil C, Albers PK, McVean G. 2019. Inferring whole-genome histories in large population datasets. Nat Genet. 51(9):1330–1338. https://doi.org/10.1038/s41588-019-0483-y
- Kendall M, Colijn C. 2016. Mapping phylogenetic trees to reveal distinct patterns of evolution. Mol Biol Evol. 33(10):2735–2743. https://doi.org/10.1093/molbev/msw124
- Le MK, Smith OS, Akbari A, Harpak A, Reich D, Narasimhan VM. 2022. 1,000 ancient genomes uncover 10,000 years of natural selection in Europe. bioRxiv. https://doi.org/10.1101/2022.08.24.505188, preprint: not peer reviewed.
- Lee KM, Coop G. 2017. Distinguishing among modes of convergent adaptation using population genomic data. Genetics. 207(4): 1591–1619. https://doi.org/10.1534/genetics.117.300417
- Lewanski AL, Grundler MC, Bradburd GS. 2024. The era of the ARG: an introduction to ancestral recombination graphs and their significance in empirical evolutionary genomics. PLoS Genet. 20(1): e1011110. https://doi.org/10.1371/journal.pgen.1011110
- Link V, Schraiber JG, Fan C, Dinh B, Mancuso N, Chiang CWK, Edge MD. 2023. Tree-based QTL mapping with expected local genetic relatedness matrices. Am J Hum Genet. 110(12):2077–2091. https://doi.org/10.1016/j.ajhg.2023.10.017
- Martin AR, Gignoux CR, Walters RK, Wojcik GL, Neale BM, Gravel S, Daly MJ, Bustamante CD, Kenny EE. 2017. Human demographic history impacts genetic risk prediction across diverse populations. Am J Hum Genet. 100(4):635–649. https://doi.org/10.1016/ j.ajhg.2017.03.004
- Maruvka YE, Kessler DA, Shnerb NM. 2011. The birth-death-mutation process: a new paradigm for fat tailed distributions. PLoS One. 6(11):e26480. https://doi.org/10.1371/journal.pone.0026480
- Mathieson I, Lazaridis I, Rohland N, Mallick S, Patterson N, Roodenberg SA, Harney E, Stewardson K, Fernandes D, Novak M, et al. 2015. Genome-wide patterns of selection in 230 ancient eurasians. Nature. 528(7583):499–503. https://doi.org/10.1038/ nature16152
- Mathieson I, Terhorst J. 2022. Direct detection of natural selection in bronze age Britain. Genome Res. 32(11-12):2057–2067. https://doi. org/10.1101/gr.276862.122
- Mirzaei S, Wu Y. 2017. RENT+: an improved method for inferring local genealogical trees from haplotypes with recombination.

Bioinformatics. 33(7):1021–1030. https://doi.org/10.1093/bio informatics/btw735

- Neale Lab. 2017. Rapid GWAS of thousands of phenotypes for 337,000 samples in the UK Biobank. https://www.nealelab.is/blog/2017/7/ 19/rapid-gwas-of-thousands-of-phenotypes-for-337000-samplesin-the-uk-biobank. [accessed 2024 Dec 14].
- Nielsen R, Vaughn AH, Deng Y. 2025. Inference and applications of ancestral recombination graphs. Nat Rev Genet. 26(1):47–58. https://doi.org/10.1038/s41576-024-00772-4
- Przeworski M, Coop G, Wall JD. 2005. The signature of positive selection on standing genetic variation. Evolution. 59(11):2312–2323. https://doi.org/10.1554/05-273.1
- Racimo F, Berg JJ, Pickrell JK. 2018. Detecting polygenic adaptation in admixture graphs. Genetics. 208(4):1565–1584. https://doi.org/10. 1534/genetics.117.300489
- Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. 2014. Genome-wide inference of ancestral recombination graphs. PLoS Genet. 10(5): e1004342. https://doi.org/10.1371/journal.pgen.1004342
- Robinson DF, Foulds LR. 1981. Comparison of phylogenetic trees. Math Biosci. 53(1-2):131–147. https://doi.org/10.1016/0025-5564(81)90043-2
- Robinson MR, Hemani G, Medina-Gomez C, Mezzavilla M, Esko T, Shakhbazov K, Powell JE, Vinkhuyzen A, Berndt SI, Gustafsson S, et al. 2015. Population genetic differentiation of height and body mass index across Europe. Nat Genet. 47(11):1357–1362. https://doi.org/10.1038/ng.3401
- Slatkin M, Rannala B. 1997. Estimating the age of alleles by use of intraallelic variability. Am J Hum Genet. 60:447–458.
- Sohail M, Maier RM, Ganna A, Bloemendal A, Martin AR, Turchin MC, Chiang CW, Hirschhorn J, Daly MJ, Patterson N, et al. 2019. Polygenic adaptation on height is overestimated due to uncorrected stratification in genome-wide association studies. Elife. 8:e39702. https://doi.org/10.7554/eLife.39702
- Speidel L, Cassidy L, Davies RW, Hellenthal G, Skoglund P, Myers SR. 2021. Inferring population histories for ancient genomes using genome-wide genealogies. Mol Biol Evol. 38(9):3497–3511. https://doi.org/10.1093/molbev/msab174
- Speidel L, Forest M, Shi S, Myers SR. 2019. A method for genome-wide genealogy estimation for thousands of samples. Nat Genet. 51(9): 1321–1329. https://doi.org/10.1038/s41588-019-0484-x
- Stern AJ, Speidel L, Zaitlen NA, Nielsen R. 2021. Disentangling selection on genetically correlated polygenic traits via whole-genome genealogies. Am J Hum Genet. 108(2):219–239. https://doi.org/10. 1016/j.ajhg.2020.12.005
- Tavaré S. 1984. Line-of-descent and genealogical processes, and their applications in population genetics models. Theor Popul Biol. 26(2):119–164. https://doi.org/10.1016/0040-5809(84)90027-3
- Uricchio LH, Kitano HC, Gusev A, Zaitlen NA. 2019. An evolutionary compass for detecting signals of polygenic selection and mutational bias. Evol Lett. 3(1):69–79. https://doi.org/10.1002/evl3.97
- Veller C, Coop GM. 2024. Interpreting population- and familybased genome-wide association studies in the presence of confounding. PLoS Biol. 22(4):1–35. https://doi.org/10.1371/journal. pbio.3002511
- Veller C, Przeworski M, Coop G. 2024. Causal interpretations of family GWAS in the presence of heterogeneous effects. Proc Natl Acad Sci U S A. 121(38):e2401379121. https://doi.org/10.1073/ pnas.2401379121
- Volz EM, Kosakovsky Pond SL, Ward MJ, Leigh Brown AJ, Frost SDW. 2009. Phylodynamics of infectious disease epidemics. Genetics. 183(4):1421–1430. https://doi.org/10.1534/genetics.109.106021
- Wakeley J. 2016. Coalescent Theory: An Introduction. Macmillan Learning.

- Wang X, Christ R, Young E, Kang CJ, Das I, Belter Jr EA, Laakso M, Aslett LJ, Steinsaltz D, Stitziel NO, et al. 2024. Genealogy based trait association with locater boosts power at loci with allelic heterogeneity. medRxiv pp. 2024–11. https://doi.org/10.1101/2024. 11.04.24316696, preprint: not peer reviewed.
- Whitehouse LS, Ray DD, Schrider DR. 2024. Tree sequences as a general-purpose tool for population genetic inference. Mol Biol Evol. 41(11):msae223. https://doi.org/10.1093/molbev/msae223
- Wohns AW, Wong Y, Jeffery B, Akbari A, Mallick S, Pinhasi R, Patterson N, Reich D, Kelleher J, McVean G. 2022. A unified genealogy of modern and ancient genomes. Science. 375(6583): eabi8264. https://doi.org/10.1126/science.abi8264
- Wong Y, Ignatieva A, Koskela J, Gorjanc G, Wohns AW, Kelleher J. 2024. A general and efficient representation of ancestral recombination graphs. Genetics. 228(1):iyae100. https://doi.org/10.1093/ genetics/iyae100

- Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, Chu AY, Estrada K, Luan J, Kutalik Z, et al. 2014. Defining the role of common variation in the genomic and biological architecture of adult human height. Nat Genet. 46(11):1173–1186. https://doi. org/10.1038/ng.3097
- Zhang BC, Biddanda A, Gunnarsson ÁF, Cooper F, Palamara PF. 2023. Biobank-scale inference of ancestral recombination graphs enables genealogical analysis of complex traits. Nat Genet. 55(5): 768–776. https://doi.org/10.1038/s41588-023-01379-x
- Zhu J, Kalantzis G, Pazokitoroudi A, Chen H, Sankararaman S, Palamara PF. 2024. Fast variance component analysis using large-scale ancestral recombination graphs. bioRxiv. pp. 2024–08 https://doi.org/10. 1101/2024.08.31.610262, preprint: not peer reviewed.

Editor: S. Gravel