

# A forest is more than its trees: haplotypes and ancestral recombination graphs

Halley Fritze,<sup>1</sup> Nathaniel Pope ,<sup>2</sup> Jerome Kelleher,<sup>3</sup> Peter Ralph <sup>1,2,4,\*</sup>

<sup>1</sup>Department of Mathematics, University of Oregon, Eugene, Oregon 97403, United States

<sup>2</sup>Institute of Evolution and Ecology, Department of Biology, University of Oregon, Eugene, Oregon 97403, United States

<sup>3</sup>Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford, Oxford OX3 7LF, United Kingdom

<sup>4</sup>Department of Data Science, University of Oregon, Eugene, Oregon 97403, United States

\*Corresponding author: Institute of Ecology and Evolution, 272 Onyx Bridge, 5289 University of Oregon Eugene, OR 97403, USA. Email: plr@uoregon.edu.

Foreshadowing haplotype-based methods of the genomics era, it is an old observation that the “junction” between two distinct haplotypes produced by recombination is inherited as a Mendelian marker. In a genealogical context, this recombination-mediated information reflects the persistence of ancestral haplotypes across local genealogical trees in which they do not represent coalescences. We show how these non-coalescing haplotypes (“locally-unary nodes”) may be inserted into ancestral recombination graphs, a compact but information-rich data structure describing the genealogical relationships among recombinant sequences. The resulting ancestral recombination graphs are smaller, faster to compute with, and the additional ancestral information that is inserted is nearly always correct where the initial ancestral recombination graph is correct. We provide efficient algorithms to infer locally-unary nodes within existing ancestral recombination graphs, and explore some consequences for ancestral recombination graphs inferred from real data. To do this, we introduce new metrics of agreement and disagreement between ancestral recombination graphs that, unlike previous methods, consider ancestral recombination graphs as describing relationships between haplotypes rather than just a collection of trees.

**Keywords:** genealogy; tree sequence; haplotypes; ancestral recombination graph

## Introduction

Ancestral recombination graphs (ARGs) describe how a set of sampled sequences are related to each other at each position of the genome in a recombining species (Brandt et al. 2024; Lewanski et al. 2024; Wong et al. 2024; Nielsen et al. 2025), and there has been significant recent progress on inference through a range of approaches (Rasmussen et al. 2014; Kelleher et al. 2019; Speidel et al. 2019; Zhang et al. 2023; Deng et al. 2024; Gunnarsson et al. 2024). One way of viewing ARGs is as a sequence of local trees, i.e. the genealogical trees that describe how each portion of the genome was inherited by the focal genomes. This is reflected in methodology of some ARG inference methods and in metrics used to assess inference accuracy, as well as in basic terminology. For instance, the “succinct tree sequence”, introduced by Kelleher et al. (2016), is a common format for describing these inferred ARGs, and is seeing wide use thanks in part to its efficiency and accompanying reliable toolkit, *tskit* (Ralph et al. 2020; Kelleher et al. 2024).

However, an ARG is emphatically not merely a sequence of trees: viewed another way, it describes inheritance relationships between ancestral haplotypes. These two points of view are related because a single haplotype may extend over many local trees; in other words, the internal nodes in the trees are labelled, and many of these labels are shared between adjacent trees (Wong et al. 2024).

Another reason we tend to focus on the trees is that much of our intuition about inference of relationships from genomic data comes from phylogenetics. Indeed, all methods might very roughly be summarized as “more similar sequences are more closely related”. For instance, two sequences that share a derived mutation are probably more closely related over some span of genome surrounding the location where the mutation occurs. It has long been observed that not only mutations but also the “junctions” between distinct haplotypes, if they could be somehow identified, would be inherited as Mendelian markers (Fisher 1954; Chapman and Thompson 2003). In more modern terminology, even in the absence of new mutations, recombination between distinct haplotypes can create a novel haplotype whose relationships and origination time could be inferred.

Haplotype identity has been largely overlooked in the literature on ARG inference—most methods that have been used so far to measure accuracy of inferred ARGs depend only on the sequence of local trees, not on how ancestral haplotypes span across these trees. For instance, Kelleher et al. (2019) and Zhang et al. (2023) compared true and inferred ARGs using average Robinson-Foulds (Robinson and Foulds 1981) and Kendall–Colijn (Kendall and Colijn 2016) distances between trees across a regular sequence of genomic positions, using sampled genotypes as labels, while Brandt et al. (2022) compared times to most recent common

Received on 29 May 2025; accepted on 12 August 2025

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

This is an Open Access article distributed under the terms of the Creative Commons Attribution-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](mailto: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](mailto:journals.permissions@oup.com).

ancestor between pairs of sampled genomes. Neither is affected by shared haplotype structure—two ARGs could be identical by either measure but imply completely different patterns of haplotype sharing and inheritance. Also, [Deng et al. \(2021\)](#) evaluated agreement of distributions of distances along the genome between tree topology changes, and [Zhang et al. \(2023\)](#) defined a generalization of Robinson-Foulds distance that is the total variation distance between the induced distribution on genotypes; however, neither of these measure the sharing of haplotypes between adjacent trees. An exception is [Ignatieva et al. \(2025\)](#), who compared distributions of haplotype spans in true and inferred ARGs, as well as more sophisticated summaries of edges. The additional information provided by haplotype structure can be important: for instance, haplotypes that extend over many local genealogies “tie together” those genealogies, allowing estimates of times of particular ancestors to be informed by larger portions of the genome on which there are many genealogies.

In this paper, we study various aspects of haplotype identity in ARGs. First, we describe a deterministic algorithm that extends the genomic region spanned by ancestral haplotypes using the principle that intermediate nodes in inheritance paths should remain unchanged when possible. These extended portions of ancestral haplotypes manifest as unary nodes in the local trees. To quantify how accurate the new information is, we define and describe how to compute new measures of (dis-)agreement between ARGs that are motivated by the Robinson-Foulds distance between trees but account for haplotype identity. These measures show that the vast majority of these extended haplotypes are correct if the trees are correct, and that substantial information about haplotypes is contained in these nodes in inferred trees as well.

## Motivation and statement of problem

Consider the (small portion of a) hypothetical ARG in [Fig. 1a](#). On the first portion of the genome (left-hand tree), the sample nodes (labelled 0, 1, and 2) coalesce into a small subtree: 1 and 2 find a common ancestor in ancestral node 3, which finds a common ancestor with node 0 in ancestral node 4. On the next portion of the genome (right-hand tree), sample node 2 has a different ancestor. This seems reasonable, and a method that infers trees separately on each portion of the genome could not be expected to produce anything different. However, the example becomes more complicated once we consider what these local genealogies imply about haplotype inheritance. [Figure 1b](#) shows the implied inheritance of haplotypes, with the haplotypes carried by 4 to the left and right of the recombination breakpoint labelled  $L$  and  $R$ . Here, sample node 2 has inherited the chunk of haplotype labelled  $L$  from ancestral node 4 via 3, and the haplotype to the right of this from some other node (and so doesn't carry haplotype  $R$ ). On the other hand, sample node 1 has inherited *both* haplotypes  $L$  and  $R$  from ancestral node 4, but the trees imply that only haplotype  $L$  is inherited via ancestral node 3. This implies—if taken literally—that there must have been a recombination event at some point between node 1 and node 4 that separated the  $L$  and  $R$  haplotypes, and then these two ancestral (and nonoverlapping) haplotypes coalesced together in ancestral node 4. Although this is possible, it seems unlikely—a more parsimonious explanation is depicted in [Fig. 1c](#), in which sample node 1 inherits the entire  $LR$  haplotype from ancestral node 4 through node 3 (and there is a recombination somewhere between node 3 and node 2). This implies that ancestral node 3 inherits from node 4 on the right-hand tree as well, which is depicted in [Fig. 1d](#)—and so node 3 has become unary in this tree. Note that the more parsimonious ARG also includes fewer edges: the three distinct edges  $4 \rightarrow 3$ ,  $3 \rightarrow 1$ , and  $4 \rightarrow 1$  in

[Fig. 1b](#) have been reduced to the two edges  $4 \rightarrow 3$  and  $3 \rightarrow 1$  in [Fig. 1d](#).

So, given the ARG shown in [Fig. 1a](#) and [b](#), it should be possible to extend the ancestral haplotype represented by node 3 to obtain the ARG shown in [Fig. 1c](#) and [d](#), thus adding additional information to the ARG. This might be surprising, as intuition from phylogenetics suggests we can only infer information about the branching points in the tree, not intermediate (unary) nodes. The goal of this paper is to answer: How can we do this, and how accurate is the resulting inference? See the Discussion for more on what this question assumes, and the connection to parsimony.

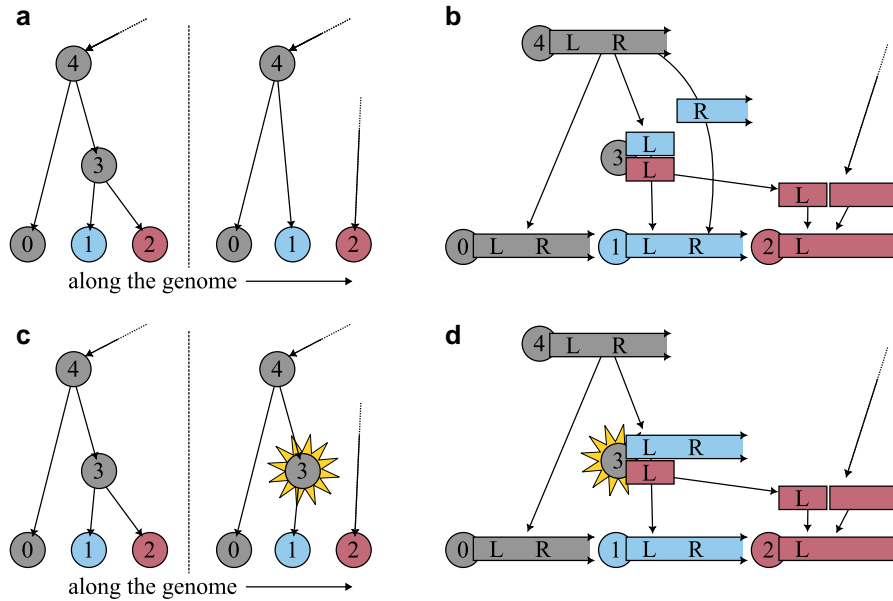
## Methods

### Notation and terminology

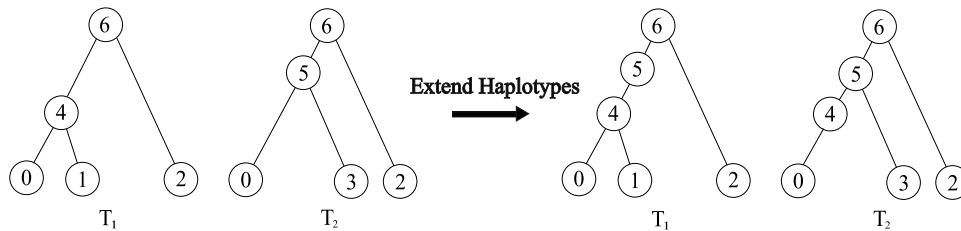
We work with the *succinct tree sequence* representation of ARGs (henceforth, “tree sequence”), to take advantage of the tools available in `tskit` ([Kelleher et al. 2024](#)), and our terminology and notation follows [Ralph et al. \(2020\)](#). For our purposes here, a tree sequence  $\mathbb{T} = (N, E)$  contains a set of nodes  $N$  which represent ancestral segments of genome, and edges  $E$  which represent relationships between nodes over different regions of the genome. Each node  $n \in N$  has a time  $t_n$ , which is the amount of time in the past that the individual who carried that segment of genome lived. Some nodes are *samples*, meaning that they represent genome sequences available as data. Each edge  $e \in E$  describes inheritance between an ancestor and a descendant over a segment of genome  $[\ell_e, r_e)$ . The ancestor and descendant are represented, respectively, by the parent node  $p_e$  and child node  $c_e$  of the edge and occur at unique times so that  $t_{p_e} > t_{c_e}$ . Suppose that the unique elements of the set of left and right edge endpoints are  $0 = a_0 < a_1 < \dots < a_n = L$ , where  $L$  is the length of the genome. Using this information, one can construct the sequence of trees  $(T_1, \dots, T_n)$  that describe how the nodes are related to each other along the genome: each  $T_k$  is a tree whose nodes are in  $N$  and that describes relationships on the half-open interval  $[a_{k-1}, a_k)$ . Nodes represent (portions of) ancestral haplotypes, and so we will use the terms interchangeably. Not all nodes appear in each tree, and we say  $n \in T_k$  for a node  $n$  if the tree  $T_k$  describes at least one parent-child relationship for node  $n$ .

### An algorithm to extend haplotypes

Given a tree sequence, our goal is to identify areas of implied inheritance of haplotypes. Generalizing from [Fig. 2](#), we do this by identifying paths of inheritance that are shared across a sequence of local trees but for which some of the intermediate nodes are missing. Concretely, suppose that if in tree  $T_k$  there is a chain of inheritance  $p \rightarrow u_1 \rightarrow \dots \rightarrow u_m \rightarrow c$  (where  $a \rightarrow b$  denotes a parent-child relationship) and in tree  $T_{k+1}$  there is a chain of inheritance  $p \rightarrow v_1 \rightarrow \dots \rightarrow v_n \rightarrow c$ , where  $\{u_i\}_{i=1}^m$  and  $\{v_j\}_{j=1}^n$  are disjoint. This situation implies that  $c$  inherited from  $p$  over the entire interval  $[a_{k-1}, a_{k+1})$ , so it seems reasonable to assume that  $c$  has inherited from  $p$  *along the same path* for that entire interval. In other words, the intermediate nodes  $\{u_i\}$  should also lie on the path from  $c$  to  $p$  in tree  $T_{k+1}$ , and conversely the nodes  $\{v_j\}$  should lie on that path in tree  $T_k$ . For instance, in [Fig. 1](#) take  $p = 4$  and  $c = 1$ , so that  $u_1 = 3$  (and  $m = 1$ ) and there are no  $v$  (so  $n = 0$ ); and so as shown in [Fig. 1b](#) and [d](#), we might extend node 3 over this entire segment. Of course, this does not always make sense—for instance, if  $u_i$  is already represented somewhere else in  $T_{k+1}$ , or if  $t_{u_i} = t_{v_j}$ , for some  $i$  and  $j$ . So, we restrict our attention to pairs of such paths in adjacent trees for which  $u_i \notin T_{k+1}$  for all  $1 \leq i \leq m$ ,  $v_j \notin T_k$  for all  $1 \leq j \leq n$ , and the times of the nodes  $\{u_i\}$  and  $\{v_j\}$  are



**Fig. 1.** A simple example showing the basic idea (described in more detail in the text): a) a small portion of an ARG without unary nodes; b) the implied inheritance pattern of the two portions of the haplotype carried by ancestral node 4, labelled L and R; c) local trees with a unary node added, which produces d) a more parsimonious haplotype inheritance pattern (that also includes fewer edges).



**Fig. 2.** A visualization of the *extend haplotypes* method. In both trees  $T_1$  and  $T_2$ , node 0 inherits from node 6, the root:  $T_1$  contains the path  $6 \rightarrow 4 \rightarrow 0$  while  $T_2$  has path  $6 \rightarrow 5 \rightarrow 0$ . The intermediate nodes 4 and 5 do not appear in  $T_2$  and  $T_1$  respectively, and so the paths are *mergeable*. The “extend haplotypes” method joins these two paths, inserting the merged path  $6 \rightarrow 5 \rightarrow 4 \rightarrow 0$  into both  $T_1$  and  $T_2$ .

unique, meaning no node in  $u_i$  has the same time as a node in  $v_j$ . Call a pair of such paths *mergeable*. So, the goal of our algorithm is to iterate over trees, identify mergeable pairs of paths, and then extend the nodes  $\{u_i\}$  to  $T_{k+1}$ . (We also extend  $\{v_j\}$  to  $T_k$ , but on a backwards pass.)

An efficient algorithm to do this is described in [Algorithm 1](#). The algorithm considers each tree transition from  $T_k$  to  $T_{k+1}$  in turn, updating its internal state (which includes possibly modifying  $T_{k+1}$ ) as it goes. Suppose we are at the transition from tree  $T_k$  to tree  $T_{k+1}$ , which is done by first removing a set of edges  $O$  and then adding another set of edges  $I$ .  $O$  defines a sub-forest  $F_O$  of  $T_k$ , and  $I$  defines a sub-forest  $F_I$  of  $T_{k+1}$ . In [Fig. 1](#) (in the portion of the tree shown), the removed edges that make up  $F_O$  are  $4 \rightarrow 3$ ,  $3 \rightarrow 1$ , and  $3 \rightarrow 2$ , and the added edges that make up  $F_I$  are  $4 \rightarrow 1$  and the new edge leading to 2. The key step in the algorithm is to determine whether the pair of paths that terminate in a given node in two adjacent trees are mergeable. The algorithm we use to do this is given as [Algorithm 2](#), and works as follows. If a pair of paths is mergeable, then the edges of the two paths must lie in  $O$  and  $I$ , respectively. Suppose an edge in  $I$  has child  $c$ . To see if  $c$  is the base of a pair of mergeable paths, the algorithm traverses up from  $c$  in both  $F_O$  and  $F_I$ ; the variables  $y_i$  and  $y_o$  indicate whether the next node upwards in the traversals (i.e. the next older node, as recall  $t_n$  is amount of time in the past that  $n$  lived) is in  $F_I$  or  $F_O$ . The traversals terminate if a node in the other tree is found (i.e. if the node traversed in  $F_O$  is in  $T_{k+1}$  or if the node traversed in  $F_I$  is in  $T_k$ )

or if a pair of traversed nodes have the same time. If these two traversals end in the same node  $p$ , the paths are mergeable. Iterating over all edges in  $O$  will thus find all mergeable pairs of paths. There is often more than one pair of mergeable paths in a tree transition; so, the algorithm merges pairs of mergeable paths, starting with pairs that add the smallest number of new edges, until no more are found.

Note that the algorithm could be applied to an undated ARG, with some adjustments. The algorithm only uses node time for two reasons: convenience, when iterating jointly up the two paths; and to avoid illegal conditions: if we tried to merge two paths on which some  $u_i$  had the same time as some  $v_j$ , then we would have a parent and child with the same time, which is not allowed.

Algorithm 1 simplifies the full algorithm implemented in software in several ways for the sake of clarity—for instance, the bookkeeping required to keep track of  $T_k$  and  $T_{k+1}$  is omitted. Furthermore, as described the algorithm does one left-to-right pass over the tree sequence; in practice we do repeated passes in both directions until no changes can be made. We require repeated passes because the additional structure added by a given path of the algorithm may introduce more mergeable paths. Most of these cases seem to occur for paths with large time differences between the parent and child nodes, or when many nodes in a path are coalescent. The number of required passes is in practice small: see empirical results below.

The main step that is omitted is a description of the Merge operation, which performs the actual extending of haplotypes. This algorithm is essentially the same as MergeNum in Algorithm 2, except with additional bookkeeping. Roughly speaking, the algorithm traverses up from the shared base node  $c$ , doing the appropriate operations to insert the nodes along the path found in  $T_k$  into the path in  $T_{k+1}$ . To do this, some edges that end at  $a_k$  will be extended to end at  $a_{k+1}$ ; some edges that begin at  $a_k$  will be postponed to begin at  $a_{k+1}$ , and some entirely new edges may be added, as in Fig. 2. Furthermore, the forests  $F_O$  and  $F_I$  (and hence  $T_{k+1}$ ) need to be updated.

**Algorithm 1:** Extend haplotypes. Given an ARG  $\mathbb{T}$  with  $N$  trees  $T_1, \dots, T_N$  for which edges  $I_k$  are added to transition from  $T_k$  to  $T_{k+1}$ , identify and merge all mergeable paths (see text). Each child node  $c_e$  of each removed edge  $e$  is checked to see if it is at the base of two mergeable paths; paths that add fewer new edges are merged first because MergeNum returns the number of new edges required to merge the two paths, and  $M$  is always less than or equal to the minimum number of new edges across all nodes in the current tree transition.

```

1 def ExtendHaplotypes( $\mathbb{T}$ ): /* For pairs of trees in the ARG,
   find mergeable paths. */
2 for  $k$  in  $1 \dots N - 1$ :
3   Set  $M = 0$  and  $M' = \infty$ .
4   while  $M < \infty$ :
5     for  $e \in I_k$ :
6       Set  $m_e = \text{MergeNum}(c_e, T_k, T_{k+1})$ . /* Find the number
   of new edges
   required. */
7     if  $m_e \leq M$ :
8       Merge( $c_e, T_k, T_{k+1}$ ) /* Extended haplotypes from
    $T_k$  into  $T_{k+1}$ . */
9     else:
10      Set  $M' = \min(m_e, M')$ .
11   Set  $M = M'$  and  $M' = \infty$ .

```

**Algorithm 2** Given a node  $c$ , trees  $T_O$  and  $T_I$ , and sub-forests  $F_O$  and  $F_I$  such that removing  $F_O$  and adding  $F_I$  turns  $T_O$  into  $T_I$ , check to see if the paths upwards from  $c$  in  $T_O$  and  $T_I$  are mergeable. If the paths are mergeable then this returns the number of new edges that would be added by extending the path from  $T_O$  to  $T_I$ ; otherwise, this returns  $\infty$ . Let  $P_O[n]$  and  $P_I[n]$  be the parents of node  $n$  in the set of edges to be removed and added, respectively (i.e. in  $F_O$  and  $F_I$ ); these are NULL if  $n$  has no parent, and we take  $t[\text{NULL}] = \infty$ . The variable  $m_e$  will record the number of new edges to be added, and  $m$  will record the number of extended haplotypes. Note that if no haplotypes can be extended ( $m = 0$ ), we return  $m_e$  as  $\infty$  so that the path is not merged in Algorithm 1.

```

1 def MergeNum( $c, T_O, T_I$ ):
2   Let  $p_i = P_I[c]$ ,  $t_i = t[p_i]$ ,  $p_o = P_O[c]$ ,  $t_o = t[p_o]$ , and  $m_e = m = 0$ .
3   while True:
4     Set  $y_i = (p_i \neq \text{NULL}) \ \& \ (p_i \notin T_O) \ \& \ (t_i < t_o)$  /* Next node is
   from  $T_I$ . */
5     and  $y_o = (p_o \neq \text{NULL}) \ \& \ (p_o \notin T_I) \ \& \ (t_o < t_i)$  /* Next node is
   from  $T_O$ . */
6     if not ( $y_i$  or  $y_o$ ):
7       break
8     if  $y_i$ :
9       if  $P_I[c] \neq p_i$ :
10        Set  $m_e = m_e + 1$ . /*  $c \rightarrow p_i$  is a new edge. */
11        Set  $c = p_i$ ,  $p_i = P_I[p_i]$ , and  $t_i = t[p_i]$ .
12      else:
13        if  $P_O[c] \neq p_o$ :
14          Set  $m_e = m_e + 1$ . /*  $c \rightarrow p_o$  is a new edge. */
15          Set  $c = p_o$ ,  $p_o = P_O[p_o]$ ,  $t_o = t[p_o]$ 
16          and  $m = m + 1$ . /*  $c \rightarrow p_o$  would be extended. */
17    if  $m = 0$  or  $p_i \neq p_o$  or  $p_i = \text{NULL}$ :
18      Set  $m_e = \infty$ . /* The paths cannot be merged. */
19    return  $m_e$ 

```

Because the algorithm needs to take multiple passes over the tree sequence in each direction, an important practical question for this algorithm is: how many passes do we need to do? The algorithm is monotone (spans of ancestral nodes only increase), so it is guaranteed to terminate in a finite number of passes, but it is also not hard to construct pathological cases that require an arbitrary number of passes. However, experimentation suggests that in practice at most five iterations are needed before the algorithm terminates. Indeed, for even large sequences Supplementary Table S1 shows that 99% of all changes to an ARG occur after the first iteration, with the algorithm always completing after four iterations.

## Dissimilarity between ARGs

If we begin with a tree sequence containing unary nodes, it is straightforward to remove the portions of each node's span on which it is unary, extend haplotypes (Algorithm 1), and quantify how much node span was correctly or incorrectly added. However, we are also interested in whether extending haplotypes improves inferred ARGs. Since we are not aware of any current methods for measuring (dis)agreement between ARGs that take into account haplotype identity, we define a measure of *matched span* to quantify this. The method (including Equations (1)–(5)) is implemented in the `tscompare` package.

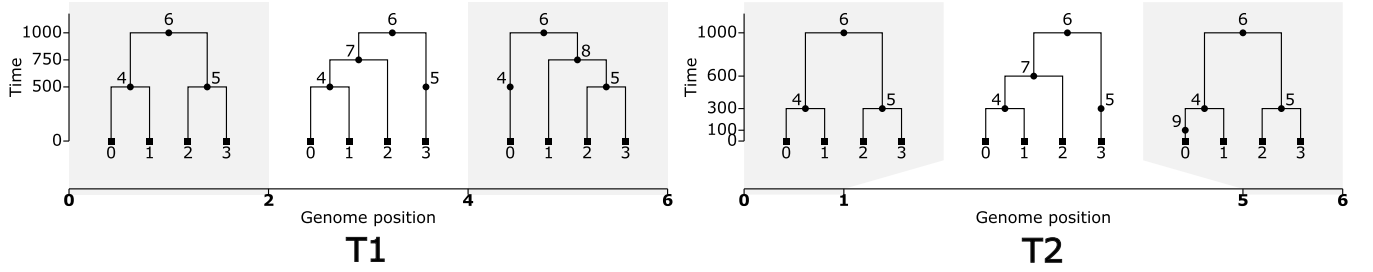
It is helpful to first describe what we compute in the simple case. We will first simulate tree sequences where nodes that are unary in local trees between coalescent haplotypes are retained. Then, each node is present in both tree sequences, and we can quantify, for each node, how much of their span is correct or incorrect by comparing to the original, true tree sequence.

Now suppose that instead of comparing two tree sequences with the same set of nodes, we wish to compare two tree sequences for which we know the sample nodes are the same but are otherwise unclear as to the equivalency of nodes across sequences. (For instance, with a simulated tree sequence and one inferred from its genotypes; nodes in the former represent actual ancestral haplotypes, and in the latter represent hypothetical ancestors which may or may not resemble the truth.) Call the two tree sequences  $\mathbb{T}_1$  and  $\mathbb{T}_2$ , which should have the same genome length and the same set of sample nodes; in what follows we think of  $\mathbb{T}_2$  as the true ARG and  $\mathbb{T}_1$  as an inferred ARG. We would like to measure (a) how much of  $\mathbb{T}_1$  is found in  $\mathbb{T}_2$ ; (b) how much of  $\mathbb{T}_2$  is found in  $\mathbb{T}_1$ ; and (c) how much of  $\mathbb{T}_1$  is *not* found in  $\mathbb{T}_2$ . (Think of these three quantities as the sizes of two relative intersections and difference between the tree sequences, thought of vaguely as sets; so note in particular that switching the roles of  $\mathbb{T}_1$  and  $\mathbb{T}_2$  gets a different set of quantities.) Roughly speaking, we first identify matching nodes as those whose sets of descendant samples agree for the largest span along the genome, and then compute for how much of their spans do their descendant samples agree (or not). An example of our method is illustrated in Fig. 3.

The method works as follows. To simplify notation suppose that the two tree sequences have the same set of breakpoints between trees, so that  $T_1^{(1)}, \dots, T_N^{(1)}$  are the trees in  $\mathbb{T}_1$  and  $T_1^{(2)}, \dots, T_N^{(2)}$  are the trees in  $\mathbb{T}_2$ . For a node  $n$  and tree  $T$  let  $S(T, n)$  denote the set of samples that inherit from  $n$  in  $T$ , and for a pair of nodes  $n_1$  and  $n_2$  with  $n_1$  in  $\mathbb{T}_1$  and  $n_2$  in  $\mathbb{T}_2$ , define

$$\mathcal{M}(n_1, n_2) = \left\{ k : S\left(T_k^{(1)}, n_1\right) = S\left(T_k^{(2)}, n_2\right) \right\},$$

to be the indices of all trees where  $n_1$  and  $n_2$  are ancestral to the same sample set in both ARGs, and



**Fig. 3.** For two tree sequences  $T_1$  and  $T_2$  the *matched span*,  $\text{match}(T_1, T_2)$ , matches nodes in  $T_1$  to nodes in  $T_2$  based on identical sample sets. We first can compute the total span of nodes in  $T_1$  and  $T_2$  as  $\|T_1\| = 46$  and  $\|T_2\| = 47$ . First, we always require that sample nodes match to their identical counterpart. Then we match intermediate nodes by comparing their descendant samples across pairs of local trees over the genome. Here we have distinct local tree pairings on the genome segments:  $[0, 1)$ ,  $[1, 2)$ ,  $[2, 4)$ ,  $[4, 5)$ , and  $[5, 6)$ . Node 4 has no match on  $[4, 5)$ , matches with node 9 on  $[5, 6)$ , and matches with node 4 in  $T_2$  everywhere else. Thus the maximal mapping for 4 should be to itself since  $m(4, 4) > m(4, 9)$ . Nodes 5, 6, and 7 all match with their identical counterpart. Lastly, node 8 has no match in  $T_2$  as there are no nodes in  $T_2$  with sample set  $\{1, 2, 3\}$ . This makes the *matched span*  $\text{match}(T_1, T_2) = 43$  and  $\text{ARF}(T_1, T_2) = \frac{3}{46}$ . Given the above matching, the *inverse matching* will match nodes 0 through 7, and node 9 has no match since its only possible match ( $4 \in T_1$ ) was not the best match from  $T_1 \rightarrow T_2$ . This means the *inverse matched span*  $\text{match}(T_1, T_2) = 43$  and  $\text{TPR}(T_1, T_2) = \frac{43}{47}$ . Additionally, let's compute  $\text{ARF}(T_2, T_1)$  and  $\text{TPR}(T_2, T_1)$ . Again, the sample nodes, root 6, and nodes 4, 5, and 7 will match to their identical counterpart and have matched nodes spans 6 (for each sample), 6, 4, 4, and 2, respectively. Node 9 matches to node 4 with a matched node span of 1. Thus,  $\text{ARF}(T_2, T_1) = \frac{41}{47}$ . The *inverse matched span* will match nodes 5, 6, 7 and sample nodes to their identical counterparts. Node 8 in  $T_1$  had no match from  $T_2$ , so it has no *inverse match*, while node 4 had two best matches from  $T_2$ : nodes 4 and 9. Since 4 in  $T_1$  shares more span with the same-numbered node in  $T_2$  than with 9, the *inverse match* for 4 in  $T_1$  is node 4 in  $T_2$ . Therefore,  $\text{TPR}(T_2, T_1) = \frac{40}{46}$ .

$$m(n_1, n_2) = \sum_{k \in \mathcal{M}(n_1, n_2)} (a_k - a_{k-1}),$$

which is the total span over which the samples below  $n_1$  in  $T_1$  matches the samples below  $n_2$  in  $T_2$ . The *matched span* of  $T_1$  in  $T_2$  is then defined to be

$$\overrightarrow{\text{match}}(T_1, T_2) = \max_{\beta: N_1 \rightarrow N_2} \sum_{n \in N_1} m(n, \beta(n)),$$

where the maximum is over all mappings  $\beta$  of nodes in  $T_1$  to nodes in  $T_2$ , and we require that samples in  $T_1$  are mapped to samples in  $T_2$ . (Note that multiple nodes in  $T_1$  may be mapped to the same node in  $T_2$ , and that some nodes in  $T_2$  may not be mapped to by any nodes in  $T_1$ .) Since the maximum is independent over nodes, we may define for each node  $n_1 \in T_1$  its *best matching node* in  $T_2$  as

$$\alpha(n_1) = \operatorname{argmax}_{n_2 \in N_2} m(n_1, n_2),$$

so that

$$\overrightarrow{\text{match}}(T_1, T_2) = \sum_{n \in N_1} m(n, \alpha(n)). \quad (1)$$

If the best-matching node is not unique, we define  $\alpha(n_1)$  to be the node in  $T_2$  out of those maximizing  $m(n_1, n_2)$  that minimizes  $|t_{n_1}^{(1)} - t_{n_2}^{(2)}|$  (and if this is not unique, we pick an arbitrary one) – however, this potential ambiguity does not affect the definition of  $\overrightarrow{\text{match}}(T_1, T_2)$ . Let  $s(T, n)$  denote the total span that node  $n$  is present in the local trees,

$$s(T, n) = \sum_{k=1}^N (a_k - a_{k-1}) \mathbf{1}_{n \in T_k},$$

where  $\mathbf{1}_{n \in T_k}$  is an indicator (i.e. it is 1 if  $n \in T_k$  and 0 otherwise), and let  $\|T_1\| = \sum_{n \in N_1} s(T_1, n)$  be the total span of all nodes in  $T_1$ . We then define the *non-matched span* of  $T_1$  in  $T_2$  by

$$\text{match}(T_1, T_2) = \sum_{n \in N_1} (s(T_1, n) - m(n, \alpha(n))) = \|T_1\| - \overrightarrow{\text{match}}(T_1, T_2),$$

which is the total span for all nodes in  $T_1$  over which their descendant samples do *not* match those of their best match in  $T_2$ .

Contrarily, given a matching  $\alpha: T_1 \rightarrow T_2$ , we want to quantify how much of  $T_2$  is represented in  $T_1$ . To do this, we define the *inverse matched span* of  $T_1$  in  $T_2$  as

$$\overleftarrow{\text{match}}(T_1, T_2) = \sum_{n_2 \in N_2} \max_{n_1 \in \alpha^{-1}(n_2)} m(n_1, n_2) \quad (2)$$

where  $\alpha^{-1}(n_2)$  is the set of all nodes  $n_1 \in T_1$  whose best match is  $n_2$ . This differs from the *matched span* of  $T_2$  in  $T_1$  because there may be more than one node in  $T_1$  that is mapped to the same node in  $T_2$  – so, if nodes  $n_1$  and  $n'_1$  are both mapped by  $\alpha$  to the same node  $n_2$ , then both count towards  $\overleftarrow{\text{match}}(T_1, T_2)$ , but only the better match counts towards  $\text{match}(T_1, T_2)$ .

A common measure of disagreement between ARGs, first proposed by [Kuhner and Yamato \(2015\)](#), is to use a weighted average Robinson-Foulds (RF) distance. This could be computed in a very similar way: instead of  $m(n, \alpha(n))$  define

$$\mathcal{M}'(n) = \left\{ k : \exists n_2 \text{ for which } s(T_k^{(1)}, n) = s(T_k^{(2)}, n_2) \right\},$$

the indices of all trees on which there is *some* node in  $T_2$  whose set of descendant samples matches those of  $n$ , and

$$m'(n, T_2) = \sum_{k \in \mathcal{M}'(n)} (a_k - a_{k-1})$$

the total span over which  $n$  finds a match. Then the average RF distance (averaged over locations in the genome) is

$$\frac{1}{L} \left( \|T_1\| + \|T_2\| - \sum_{n_1 \in N_1} m'(n_1, T_2) - \sum_{n_2 \in N_2} m'(n_2, T_1) \right).$$

(This holds because the average number of nodes per local tree in  $T_1$  is  $\|T_1\|/L$ , so the difference gives the average number of non-matching nodes.) In other words, we require a node in  $T_1$  to match the *same* node in  $T_2$  across all trees, but average RF distance allows a different node to match on each tree. The other differences are that *average* RF distance normalizes by sequence length, and is symmetrized. The RF distance between two trees was defined by [Robinson and Foulds \(1981\)](#) to be the minimum number of branch

contraction/expansion operations needed to move from one tree to the other (which they then show is equal to the number of branches that induce different splits on the labels). A similar metric on ARGs could be defined using the subgraph-prune-and-regraft moves used by [Deng et al. \(2024\)](#).

The matched span,  $\text{match}(\mathbb{T}_1, \mathbb{T}_2)$ , measures agreement between topologies, but not times. If the ARG is dated (e.g. as in [Wohns et al. 2022](#); [Deng et al. 2024](#)), we can naturally use the “best match”  $\alpha$  to also compare times. Empirically, dating error seems to be more or less homoskedastic on a log scale, so we recommend using the weighted root-mean-squared error of  $\log(\text{times})$ , computed as

$$\begin{aligned} \text{wRMSE}_t(\mathbb{T}_1, \mathbb{T}_2) &= \sqrt{\frac{\sum_{n \in N_1} s(\mathbb{T}_1, n) \left( \log(1 + t_n^{(1)}) - \log(1 + t_{\alpha(n)}^{(2)}) \right)^2}{\|\mathbb{T}_1\|}}, \end{aligned} \quad (3)$$

where the transformation is  $t \mapsto \log(1 + t)$  to avoid  $\log(0)$ . The mean is computed weighting by node span, so that a dating error is more impactful for a node with a longer span.

The implementation of this method in `tscompare` additionally produces relative values: the ARG RF value (non-matched span relative to  $\mathbb{T}_1$ ) and true proportion represented (inverse matched span relative to  $\mathbb{T}_2$ ). The ARG RF (which we call “ARF”) is defined to be the non-matched span proportional to the total span of nodes in  $\mathbb{T}_1$ ,

$$\text{ARF}(\mathbb{T}_1, \mathbb{T}_2) = 1 - \frac{\text{match}(\mathbb{T}_1, \mathbb{T}_2)}{\|\mathbb{T}_1\|}, \quad (4)$$

and so if  $\mathbb{T}_2$  represents the truth, is analogous to a false positive rate. The *true proportion represented* (TPR) is the inverse matched span between two trees relative to the total span of nodes in  $\mathbb{T}_2$ ,

$$\text{TPR}(\mathbb{T}_1, \mathbb{T}_2) = \frac{\text{match}(\mathbb{T}_1, \mathbb{T}_2)}{\|\mathbb{T}_2\|}, \quad (5)$$

and is analogous to statistical power. The outputs framed as proportions relative to one of the given tree sequences is more easily understood for comparing pairs of tree sequences than the original matched span and inverse matched span, whose units are length of spans. We give an example computing the ARF and TPR between two small tree sequences in [Fig. 3](#).

**Metrics on ARGs.** Neither the matched span or non-matched span of  $\mathbb{T}_1$  in  $\mathbb{T}_2$  are metrics in the mathematical sense (i.e. symmetric, nonzero distance between distinct points, and satisfying the triangle inequality). This is by design: in practice it is not possible to infer all aspects of the true ancestry of a set of samples (i.e. all their genetic ancestors who ever lived), and so we wanted to quantify “How much of the true relationships does this ARG represent?” Indeed, even a symmetrized version of non-matched span is not a metric: for instance, if  $\mathbb{T}_2$  is produced from  $\mathbb{T}_1$  by adding a new, entirely unary node, then  $\text{ARF}(\mathbb{T}_1, \mathbb{T}_2) = \text{ARF}(\mathbb{T}_2, \mathbb{T}_1) = 0$ . This is again by design: in general, edges in ARGs represent transmission of genetic material through many ancestors, and so addition of entirely unary ancestors on an edge is (arguably) not wrong. (Note that our metrics do distinguish these two ARGs: in this example  $\text{TPR}(\mathbb{T}_1, \mathbb{T}_2) < 1$ .)

A related (but not identical) way to construct a metric on (undated) ARGs modifies the definitions above so that  $\alpha$  is a bijection. To do this,

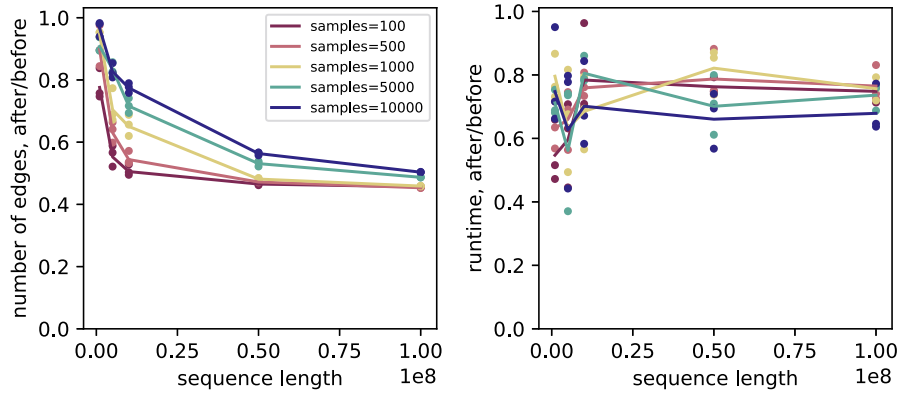
write  $N_1$  and  $N_2$  for the nodes of  $\mathbb{T}_1$  and  $\mathbb{T}_2$  respectively, and suppose that  $\varphi$  is a bijection from  $N_1$  to  $N_2$ . Since we want to allow nodes to remain unpaired, for notational convenience suppose that both  $N_1$  and  $N_2$  have appended an arbitrary number of “unrelated” nodes (since these are not ancestral to any samples they do not affect the metrics). Then, define  $d_\varphi(\mathbb{T}_1, \mathbb{T}_2)$  to be the total span of nodes  $n_1 \in N_1$  on which  $n_1$  and  $\varphi(n_1)$  do not match (i.e. are not ancestral to the same set of samples), plus the total span for nodes  $n_2 \in N_2$  on which  $n_2$  and  $\varphi^{-1}(n_2)$  do not match. The quantity  $\min_\varphi d_\varphi(\mathbb{T}_1, \mathbb{T}_2)$ , where the minimum is over bijections that preserve the partial ordering induced by the ARG, defines a metric on undated ARGs. (Actually, as written it defines a metric on the equivalence classes of undated ARGs that are obtained by removing any structure not ancestral to any samples.) To see this, note that relative ordering and descendant sample sets uniquely determine the (equivalence class of an) ARG, and that the triangle inequality is satisfied because  $d_{\varphi_{12}}(\mathbb{T}_1, \mathbb{T}_2) + d_{\varphi_{23}}(\mathbb{T}_2, \mathbb{T}_3) \geq d_{\varphi_{23 \circ \varphi_{12}}}(\mathbb{T}_1, \mathbb{T}_3)$ .

This metric,  $\min_\varphi d_\varphi(\mathbb{T}_1, \mathbb{T}_2)$ , differs from the definitions above because we require a strict one-to-one match between nodes in  $\mathbb{T}_1$  and  $\mathbb{T}_2$  (i.e.  $\varphi$  is bijective, while there may be more than one node in  $\mathbb{T}_1$  that  $\alpha$  matches to the same node in  $\mathbb{T}_2$ ). We did not require this property when defining non-matched span above for two reasons: first, the arguable non-wrongness of additional unary nodes, and second, requiring a bijection makes computation much more difficult.

The RF distance ([Robinson and Foulds 1981](#)) essentially counts the number of differing branches between two trees; the averaged RF distance ([Kuhner and Yamato 2015](#)) averages this distance across local trees, weighted by span along the genome. The method we present here for measuring dissimilarity between topologies of ARGs is a straightforward generalization that takes into account span along the genome of inferred ancestral haplotypes (and separates the metric into two pieces). However, the RF metric has many undesirable properties – for instance, moving a single tip can result in a tree with maximum distance to the original – and there is a substantial literature giving more robust generalizations (reviewed by [Labrés et al. 2021](#)). Many of these generalizations (e.g. [Böcker et al. 2013](#)) relax the requirement that the match between subtended sample sets be exact, and weight matches in some way by the size of the dissimilarity. We considered such definitions as well, but kept to the simple case for computational tractability – the generalization of [Böcker et al. \(2013\)](#) is NP-hard to compute, even for a single tree. In the ARG literature, [Zhang et al. \(2023\)](#) defines a metric (called “ARG total variation distance”) that includes branch lengths, in a way similar to [Robinson and Foulds \(1979\)](#) and [Kuhner and Felsenstein \(1994\)](#); however, it is still applied to ARGs as an average over local trees, without enforcement of identity across haplotypes; it would be useful to extend our dissimilarity to include branch lengths.

## Simulations

Our method for extending haplotypes is applicable to any ARG. However its accuracy depends on the overall structure of the ARG it is applied to. Thus to understand how well our methods can infer ancestral haplotypes we work with ARGs simulated across a range of parameter values. To do this, we simulate ARGs containing full haplotypes using Hudson’s algorithm as implemented in `msprime` ([Kelleher et al. 2016](#); [Baumdicker et al. 2021](#)), with the `coalescing_segments_only` option set to `False`. Although `msprime` simulates many events that do not create a coalescence in some local tree, by default it only outputs information for nodes which contain a coalescence (i.e. are the MRCA of some pair of samples at some point on the genome).



**Fig. 4.** Ratio of (left) number of edges, and (right) runtime for computing Tajima’s  $D$ ; before and after extending haplotypes. For instance, extending haplotypes reduces number of edges by about 50% and statistic computation runtime by about 20% for long sequences. Horizontal axis shows sequence length; colours show numbers of samples; with lines showing averages across replicates. The original tree sequence was simulated with the “expanding dog” expanding population and subset to various sizes; see Methods for details. Absolute values are shown in [Supplementary Fig. S1](#).

Furthermore, by default it only outputs those segments of the genome on which there is a coalescence. Said another way, by default all ancestral nodes in an ARG output by `msprime` are the MRCA of some pair of samples at every point in the genome on which they are represented. However, here we are interested in those segments of genome on which the nodes are *not* coalescent; i.e. where they are unary in the local trees. Setting `coalescing_segments_only` to `False` includes just this information: any ancestral segments for which these coalescent nodes are ancestral to any samples – so, the unary portions of their spans as well. However, this includes more information than we want: we hope to recover those portions of ancestral haplotypes on which the nodes are unary, but adjacent to a region of the genome where the node is not unary. For instance, if a lineage carrying an ancestral segment of genome that spans  $[a, c]$  coalesces with another spanning  $[b, d]$ , with  $a < b < c < d$ , then the resulting node is only coalescent on  $[b, c]$  but we hope using this algorithm to extend the node’s span to  $[a, b]$  and  $[c, d]$  (on which the node is unary). However, following this example, the first lineage might also carry a segment  $[x, y]$  that is disjoint from the segment  $[a, d]$ . We call these segments “isolated non-coalescent segments”; they have also been called “trapped unary spans” (by [Wong et al. 2024](#)). Such isolated segments will not be recovered by our algorithm, and would likely be unrecoverable by any other method. So, after simulation, we first remove these isolated, non-coalescent segments. To give an idea of what proportion of the full spans of ancestral nodes these isolated non-coalescent segments represent, a simulation of 1,000 samples with genome length  $5 \times 10^7$ , recombination rate  $10^{-8}$ , and population size  $10^4$  has about half the total span of all nodes in isolated, non-coalescent segments. For more discussion of these segments, see [Baumdicker et al. \(2021\)](#).

We used simulations of several scenarios. To include the effects of heterogeneous recombination rate, in some we used `stdpopsim` ([Adrion et al. 2020](#)) to simulate chromosome 1 of *Canis familiaris* using the CanFam3 genetic map from [Campbell et al. \(2016\)](#). “Constant dog” simulations simulated this chromosome in a population of (constant) size  $10^4$ . “Expanding dog” simulations were similar, but used a discrete-time Wright–Fisher model to simulate a small population of 100 that expanded to 1,000 individuals ten generations ago, which then doubled every generation to reach 512,000 individuals. All jobs for which runtime was recorded were executed on an Intel Xeon Gold 6,148 processor. Additionally, using our matched span methods,

we compare accuracy between sequences modified from a “true” ARG, (see [Fig. 6](#)). These ARGs were simulated with an effective population size of 10,000 and recombination rate of  $10^{-8}$  with between 10 and 1,000 samples and a genome length between  $10^6$  to  $5 \times 10^7$ .

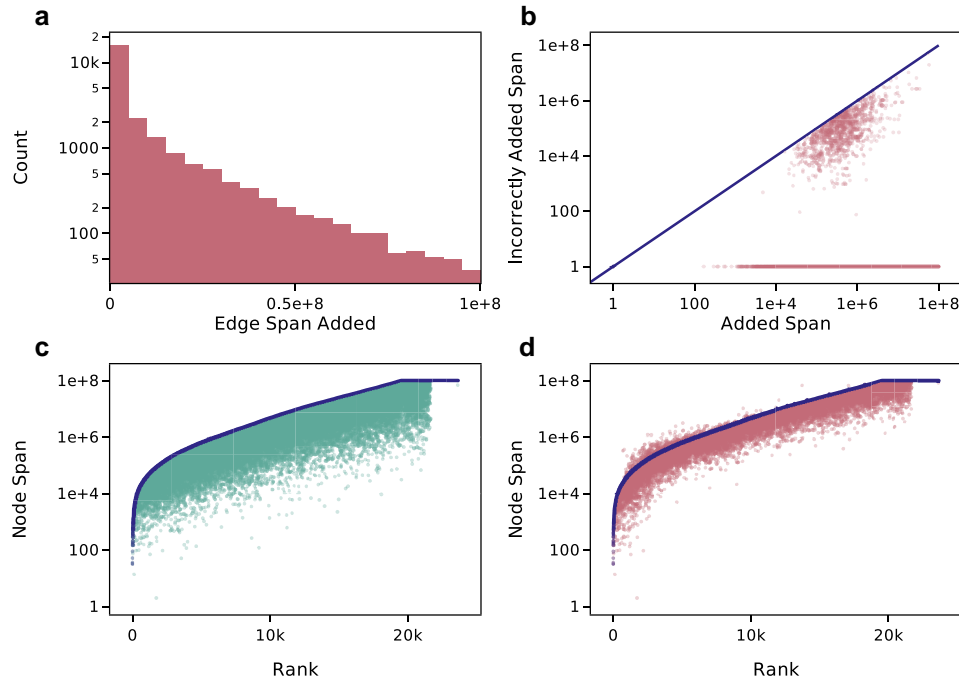
## Results

We next present the results of various experiments with inferred ARGs. First, we quantify the reduction in number of edges and resulting computational speedups. Next, to validate the principle of the algorithm, we describe accuracy of inferred additional edge spans when applied to true ARGs with unary spans removed. Finally, we explore implications for accuracy with inferred ARGs.

### Tree sequence compression and computation

In the simple example in [Fig. 1](#), extending haplotypes replaces three edges ( $0 \rightarrow 3$  and  $3 \rightarrow 4$  on the left tree, and  $0 \rightarrow 4$  on the right tree) by two edges ( $0 \rightarrow 3$  and  $3 \rightarrow 4$  on both trees). If all edge endpoints were unique, then we’d expect *every* edge to be extendable on one of its ends (except those pendant to the root and some of those adjacent to chromosome ends), leading to a reduction in number of edges by almost exactly one third. Experiments with an earlier version of the algorithm showed that if we only extend haplotypes on such “paths of length 1”, then the hypothesized reduction of  $1/3$  is achieved for long sequences. It is possible for extending haplotypes with [Algorithm 1](#) to add edges, as in [Fig. 2](#), but we still expect the number of edges to decrease by more than  $1/3$ . Indeed, [Fig. 4](#) shows that [Algorithm 1](#) nearly cuts the number of edges in half, as long as the sequence is long enough.

This reduction in edges can also lead to a reduction in computation time for algorithms using the succinct tree sequence data structure. Indeed, [Fig. 4](#) shows that computation time is reduced by 10–20% for a typical statistic (here, Tajima’s  $D$ ), computed in an efficient incremental manner along the genome as implemented in `tskit`. As described in [Ralph et al. \(2020\)](#), for these incremental algorithms the addition or removal of an edge requires updates to the state of the parent node and all nodes ancestral to it. Extending haplotypes yields a tree sequence with fewer edge removals and insertions, and thus requires fewer traversals to the roots.



**Fig. 5.** The effect of extending haplotypes on per-node spans in an ARG simulated with  $10^4$  diploid samples in a population with  $N_e = 10^4$ , and recombination rate of  $10^{-8}$  on a sequence of length  $10^8$ . a) Distribution of total amount of span added across nodes by extending haplotypes with [Algorithm 1](#); note the log scale on the y axis. b) Amount of incorrectly added span, plotted against total span, by node. 95% of nodes have no incorrect span; of the remainder, nearly all have less than 5% incorrectly added; see [Supplementary Fig. S4](#). Note the log scale on both the x and y axes. Plots c) and d) show total spans per node, ordered by total span in the original ARG (which includes unary nodes). Dark dots forming a line in both show the span of each node in this original ARG. In c), lighter scattered dots show node span after removing unary spans using *simplify*, while in d), lighter scattered dots show node span after extending haplotypes, i.e. applying [Algorithm 1](#) to the simplified ARG.

[Supplementary Fig. S2](#) shows these results are not specific to the demographic scenario. [Supplementary Figs. S1](#) and [S3](#) also show that our implementation of [Algorithm 1](#) is quite efficient, running at chromosome scale in seconds to minutes for hundreds or thousands of samples, or minutes to hours for tens of thousands of samples.

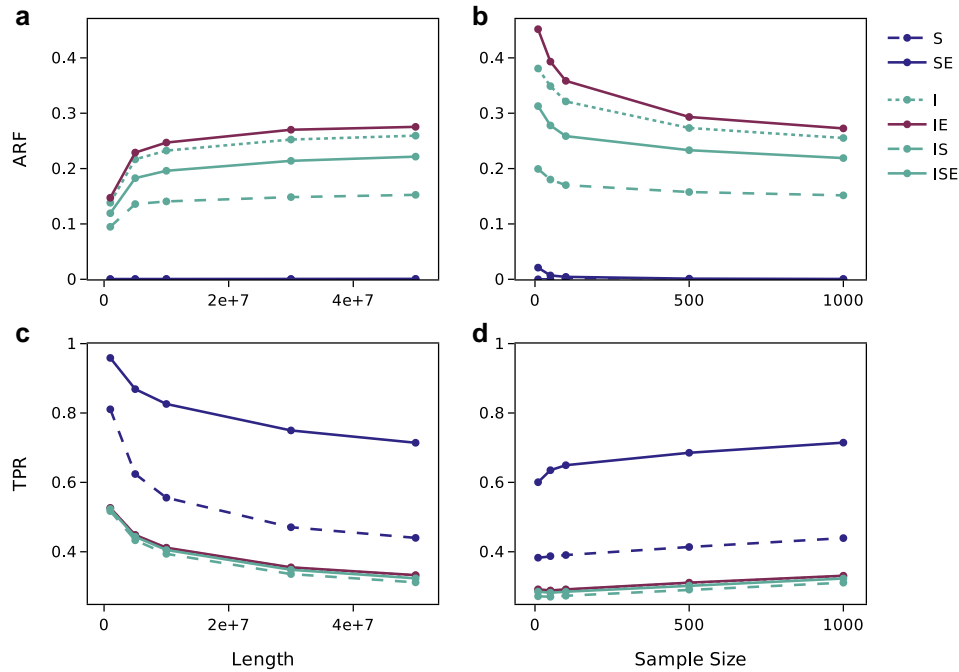
### Accuracy with true trees

Our next task is to confirm that the haplotypes extended by [Algorithm 1](#) are indeed correct – i.e. that in addition to compression, we are also gaining information. To do this, we simulate ARGs containing full haplotypes using Hudson’s algorithm in *msprime*, apply the simplification algorithm ([Kelleher et al. 2018](#); [Wong et al. 2024](#)) to reduce these so that there are no unary nodes (i.e. any node present in a local tree is a coalescent node or a sample), and then apply [Algorithm 1](#) to the result (see Methods for more detail). The method can potentially extend the spans of each node (additional span over which the node will be unary); and we can quantify how much of these extended spans were in the original ARG (and thus correctly extended).

As seen in [Fig. 5](#), the vast majority of span added by extending haplotypes is correct. In this example (which is typical), 99% of all added span is correct; 95% of nodes have no incorrectly added span; and those incorrectly added spans are nearly always a small fraction of the original span ([Fig. 5b](#) shows that incorrectly added span is often a factor of 10–100 smaller than the added span). The added information is significant: the algorithm typically increases spans (i.e. lengths of ancestral haplotypes) by around 50%. For instance, [Fig. 5a](#) shows that tens of megabases have been added to the spans of many nodes; while comparing [Figs. 5c](#) to [5d](#) shows that much of the span removed by simplification has been

replaced – in those figures, each point is one node; span after extending (points in D) are much closer to the original full span (solid blue points forming a line) than is span after simplifying (points in C). Some of the incorrectly added spans may be because of violations of the parsimony assumption, while others are due to situations with ambiguous information (and hence more than one way to extend a haplotype). It may be possible to improve the behaviour of the algorithm in the latter situation, at the cost of a more complex algorithm.

These statistics are also reflected at the genomic scale, using measures of matched span. Line “SE” in [Fig. 6](#) shows total amounts of span removed by simplification and re-inferred by extending haplotypes with [Algorithm 1](#) (correctly and incorrectly). (Lines labelled with “I” involve re-inference of the ARG; discussed next.) The top row shows the proportion of the given ARG that does not match the original (“ARF”), showing that the total amount of mis-matching span produced by extending haplotypes is very small ( $\approx 1\%$ ). (Simplification does not produce non-matched span, so the “S” line is at zero.) The bottom row shows the proportion of the original ARG that is represented in the given ARG (“TPR”). This shows that a large proportion of an ARG can be removed by simplification, indicating that coalescent nodes are unary over a substantial portion of their spans. (Since the *simplify* operations removes these unary portions of haplotypes, the “simplified” line (“S”) on the bottom two plots of [Fig. 6](#) shows the proportion of nodes’ spans on which they are not unary.) However, [Algorithm 1](#) can correctly replace most of these portions of haplotypes, especially with larger sample sizes and sequence lengths (the “simplified–extended” line; “SE”). For instance, the rightmost points show that with 1,000 samples and a  $5 \times 10^7$ bp genome, haplotypes are unary on about half their spans (on average), and



**Fig. 6.** Accuracy and sensitivity of extended haplotypes across a range of sample sizes, on a sequence of length  $5 \times 10^7$  (right; b, d); and a range of sequence lengths, with 1,000 diploid samples (left; a, c). For each, a simulated ARG containing unary haplotype spans was (i) *simplified* (“S”), removing the unary spans, and (ii) *inferred* (“I”), using `tsinfer` on genotypes and dated using `tsdate`; then each of these had its haplotypes *extended* (“SE”, “IE”). The inferred, then simplified, ARG (“IS”) and its subsequent extension (“ISE”) are also shown. **Top row** (a,b): ARF, the dissimilarity to the true ARG, as proportion of haplotypes that are not represented in the true ARG (Equation (4)). **Bottom row** (c,d): TPR, agreement to the true ARG as proportion of the true ARG that is represented.

extending haplotypes can infer more than half of this missing unary span from coalescent information only.

While our results are accurate at both the individual and genomic scale, we suspect that corrections can be made to further increase accuracy. [Algorithm 1](#) relies on only the structure of the ARG as a tree sequence (its nodes, edges, and local trees), and does not use any other parameters for inference. However, it is possible that information including the recombination rate or gene conversion rate of an ARG could improve accuracy and reduce the amount of incorrectly added span for individuals. We decided to omit incorporating this information for simplicity and efficiency of our construction.

## Inferred ARGs

So far, we have demonstrated that there is potentially ample information in the coalescent-only trees to extend haplotypes. Does this work with *inferred* ARGs? As illustrated by [Wong et al. \(2024\)](#), there is a significant diversity in the structures inferred by current methods, and here we focus on `tsinfer` [Kelleher et al. \(2019\)](#) which infers ARGs containing unary nodes as a byproduct of its inference algorithm. A comparison across other ARG inference methods is left for future work. We simulated ARGs containing unary-spanning haplotypes (as above), then re-inferred ARGs from the associated genotypes and performed various operations on the results. Since our matched spans method breaks ties with time, we date the re-inferred ARG using `tsdate`. First, [Fig. 6](#) shows that `tsinfer` has a substantial portion of already-extended haplotypes: comparing the “inferred” (“I”; dotted green) line to the “inferred-simplified” (“IS”; dashed green) line we see that inclusion of these unary spans increases the amount of correctly inferred material by around 4% (bottom panels), but that roughly 20% of the unary spans in the inferred ARG are incorrect (top panels). Furthermore, comparing to the inferred-extend (“IE”; red) line, we

see that extending the tree sequence output with `tsinfer` adds relatively little span (less than about 1%). Extending haplotypes (with [Algorithm 1](#)) of the inferred-and-then-simplified ARG (“ISE”; solid green line) produces an ARG with both less correct and incorrect span. Additionally, the “IE” and “ISE” ARGs contain approximately the same number of edges (difference of  $\approx 1\%$ ). It is also helpful to note that accuracy (i.e. proportion of the true ARG that is inferred; bottom panels) greatly increases with larger sample sizes, possibly due to resolution of polytomies. (Recall that due to computational constraints, “correct” and “incorrect” spans are determined here by exact match of subtending samples.) We additionally provide values of “I”, “IS”, “IE”, and “ISE” in [Supplementary Table S2](#), and numbers of edges and runtime for Tajima’s D in [Supplementary Fig. S5](#).

## Discussion

We began this study with the observation that the simple transformation of [Fig. 2](#) would reduce the number of edges in the succinct tree sequence representation ARGs. This is essentially a recombination-based parsimony argument, and we have shown that this line of reasoning leads to ARGs that are substantially more compact and faster to operate on, and that contain more complete information about true ancestral relationships. These extended ancestral haplotypes manifest as unary nodes in the local trees. Although a number of ARG inference methods may be taking advantage of this information, it is our impression that this source of information is not widely appreciated. In fact, due to the field’s focus on local trees rather than haplotypes, we had to develop a haplotype-aware measure of (dis)agreement between ARGs in order to study the accuracy of the proposed algorithm.

There are good reasons to think that lengthening the spans of ancestral haplotypes could lead to substantial gains in accuracy of ARG inference. For instance, information about inferring the age of a

particular mutation derives almost entirely from constraints at nearby, linked sites. Extending ancestral haplotypes from one site into neighbouring regions conceptually allows information from those local trees to inform age inference at that site as well.

We have also explored the degree to which `tsinfer` already makes use of this information, and whether this algorithm can be used to improve inference. The results do not provide a clear ordering: for instance, although `tsinfer`-produced ARGs have a substantial portion of correctly inferred unary haplotypes, removing these with simplification decreases both ARF (i.e. proportion of “wrong” haplotypes) and TPR (the proportion of the truth that is correctly inferred). Extending haplotypes restores a large amount of this correctly inferred span, but also introduces incorrect spans. Presumably, this occurs because both correctly and incorrectly inferred haplotypes are extended. Further work is needed to determine how the balance of “true and false positive rates” affects downstream uses, and whether results would differ if the requirement that sets of subtended nodes match exactly was relaxed. The efficient computational tools we have implemented (in `tskit` and `tscompare`) should facilitate this exploration.

Another consideration is storage and computational efficiency: extending haplotypes reduces the number of edges, and thus file-size and (usually) runtime for computation. Figure 4 shows that the case is clear for simulated ARGs, substantially reducing both size and runtime. However, Supplementary Fig. S5 shows the situation is more complex for inferred ARGs: extending haplotypes on a `tsinfer`-produced ARG actually increases runtime somewhat, but simplifying-then-extending reduces runtime dramatically while keeping the number of edges roughly the same. More work is needed to understand how generalizable this is and what the source of these effects are.

**Ignorance and omission in an ARG.** As motivation, we presented above a “historical” view of ARGs – i.e. that each aspect of an inferred ARG is intended to represent a portion of some particular historical genome (for instance, the MRCA of some set of sampled genomes). Furthermore, Fig. 1a and b implicitly takes the position that relationships not depicted in an ARG are implied to not exist. As discussed in Wong et al. (2024), an alternative interpretation of the ARG depicted in Fig. 1a and b would be that we have no information as to how node 2 inherited from node 4 on the right-hand interval, rather than saying that the line of transmission specifically did not pass through node 3. The “simplification” algorithm (Kelleher et al. 2018) and the Hudson algorithm for coalescent simulation (Hudson 1983; Kelleher et al. 2016) each specifically discard information about any such “non-coalescent” portions of ancestral haplotypes; so for ARGs produced by these algorithms, the correct interpretation is that the omission of unary spans reflects a lack of knowledge. In this paper, we have shown that, for the most part, this missing information can be imputed.

**Parsimony.** Much of the early work on ARG inference aimed to extract as much information as possible out of the small datasets of the time, and so, roughly speaking, integrated over possible ARGs with the goal of inferring higher-level parameters: mostly, scaled mutation rate and recombination rate (for instance, Hudson and Kaplan 1985; Griffiths and Marjoram 1996; Kuhner et al. 2000; Stephens and Donnelly 2000; Fearnhead and Donnelly 2001). However, the space of possible ARGs for a given dataset is extremely large, and other work aimed to identify the minimum number of recombinations needed to explain a given dataset under the infinite alleles model of mutation (e.g. Hein 1990; Myers and Griffiths 2003; Song and Hein 2005), which turns out to be NP-complete

(Wang et al. 2001). So, the field turned to more heuristic methods – for instance, Minichiello and Durbin (2006) used an algorithm to produce “plausible” ARGs (i.e. those that explained the data with few mutations and recombinations), and searched for associations with traits in the resulting ensemble of ARGs. (See Wong et al. (2024) for more historical discussion.) Our approach for extending haplotypes follows the same logic, that an ARG with fewer recombination events is more parsimonious, and thus more likely. For this reason, it will occasionally be wrong even if the trees are correct, although in practice this source of error is likely much smaller than error in tree inference itself.

**IBD in ARGs.** The term “identity by descent” (IBD) is used to mean many different (but related) things, and length distributions of shared IBD segments can be used for inference of recent demographic history (for instance, Browning and Browning 2015; Yang et al. 2016; Ringbauer et al. 2017; Al-Asadi et al. 2019; Silcocks et al. 2023). A commonly-used definition in the context of a given ARG says that the two genomes share an IBD segment if each has inherited the segment from their common ancestor along a single path (e.g. Ralph and Coop 2013). Largely for computational reasons, this is the definition that is used in `tskit`’s IBD-finding methods (by G. Tsambos in `tskit` Kelleher et al. 2024; see also Tsambos et al. 2023). However, many simulation and inference methods produce ARGs as shown in Fig. 1a and b, in which inheritance of a single segment is represented by more than one edge. This means that the `ibd_segments` method of `tskit` will return shorter segments than it ought to. However, as we have shown above, our method of extending haplotypes will modify the tree sequence so that the inherited segment is represented by a single edge (as in Fig. 1c and d). So, if our method (`extend_haplotypes`) is applied before finding IBD segments (with `ibd_segments`), then the resulting segments should much more accurately represent the IBD segments (in the “path” sense used here) implied by the tree sequence. Whether the resulting segments better match those predicted by theory depends also on the quality of tree inference. Note also that Huang et al. (2024) and Guo et al. (2024) both provide methods for tree sequences to compute a different definition of IBD (segments on which the MRCA does not change), which is unaffected by this issue. Further work is needed to understand how accurately IBD segments are inferred by various ARG inference methods.

## Data availability

The method to extend haplotypes described here is available through the `tskit` python and C APIs (<https://tskit.dev/tskit>) as `extend_haplotypes`; methods to compare ARGs are implemented in the `tscompare` python package (<https://tskit.dev/tscompare>). Scripts used to produce the results in this paper are available at <https://github.com/hfr1tz3/haplotypes-and-ancestral-recombination-graphs>.

Supplemental material available at GENETICS online.

## Acknowledgments

The authors would like to thank Dr. Yan Wong for his helpful comments and suggestions, as well as three anonymous reviewers for their thoughtful suggestions.

## Funding

J.K. acknowledges the Engineering and Physical Sciences Research Council (EPSRC research grant EP/X024881/1) and the Robertson

Foundation. J.K. and P.R. acknowledge support from the National Institutes of Health: National Human Genome Research Institute (NIH NHGRI research grant HG012473). The authors have no competing interests.

## Conflict of interests

None declared.

## Literature cited

- Adrion JR et al. 2020. A community-maintained standard library of population genetic models. *Elife*. 9:e54967. <https://doi.org/10.7554/eLife.54967>.
- Al-Asadi H, Petkova D, Stephens M, Novembre J. 2019. Estimating recent migration and population-size surfaces. *PLoS Genet*. 15: 1–21. <https://doi.org/10.1371/journal.pgen.1007908>.
- Baumdicker F et al. 2021. Efficient ancestry and mutation simulation with msprime 1.0. *Genetics*. 220:iyab229. <https://doi.org/10.1093/genetics/iyab229>.
- Böcker S, Canzar S, Klau GW. 2013. The generalized Robinson-Foulds metric. In Darling A, Stoye J, editors. *Algorithms in bioinformatics*. Berlin, Heidelberg. Springer Berlin Heidelberg. p. 156–169.
- Brandt DY, Huber CD, Chiang CW, Vecchyo DO-D. 2024. The promise of inferring the past using the ancestral recombination graph. *Genome Biol Evol*. 16:evae005. <https://doi.org/10.1093/gbe/evae005>.
- Brandt D, Wei X, Deng Y, Vaughn AH, Nielsen R. 2022. Evaluation of methods for estimating coalescence times using ancestral recombination graphs. *Genetics*. 221:iyac044. <https://doi.org/10.1093/genetics/iyac044>.
- Browning SR, Browning BL. 2015. Accurate non-parametric estimation of recent effective population size from segments of identity by descent. *Am J Hum Genet*. 97:404–418. <https://doi.org/10.1016/j.ajhg.2015.07.012>.
- Campbell CL, Bhérier C, Morrow BE, Boyko AR, Auton A. 2016. A pedigree-based map of recombination in the domestic dog genome. *G3 (Bethesda)*. 6:3517–3524. <https://doi.org/10.1534/g3.116.034678>.
- Chapman NH, Thompson EA. 2003. A model for the length of tracts of identity by descent in finite random mating populations. *Theor Popul Biol*. 64:141–150. [https://doi.org/10.1016/S0040-5809\(03\)00071-6](https://doi.org/10.1016/S0040-5809(03)00071-6).
- Deng Y, Nielsen R, Song YS. 2024. Robust and accurate Bayesian inference of genome-wide genealogies for large samples. *Nat Genet*. 57:2124–2135.
- Deng Y, Song YS, Nielsen R. 2021. The distribution of waiting distances in ancestral recombination graphs. *Theor Popul Biol*. 141:34–43. <https://doi.org/10.1016/j.tpb.2021.06.003>.
- Fearnhead P, Donnelly P. 2001. Estimating recombination rates from population genetic data. *Genetics*. 159:1299–1318. <https://doi.org/10.1093/genetics/159.3.1299>.
- Fisher RA. 1954. A fuller theory of ‘junctions’ in inbreeding. *Heredity (Edinb)*. 8:187–197. <https://doi.org/10.1038/hdy.1954.17>.
- Griffiths R, Marjoram P. 1996. Ancestral inference from samples of DNA sequences with recombination. *J Comput Biol*. 3:479–502. PMID: 9018600. <https://doi.org/10.1089/cmb.1996.3.479>.
- Gunnarsson ÁF et al. 2024. A scalable approach for genome-wide inference of ancestral recombination graphs [preprint]. *bioRxiv*. <https://doi.org/10.1101/2024.08.31.610248>.
- Guo B et al. 2024. Strong positive selection biases identity-by-descent-based inferences of recent demography and population structure in *plasmodium falciparum*. *Nat Commun*. 15:2499. <https://doi.org/10.1038/s41467-024-46659-0>.
- Hein J. 1990. Reconstructing evolution of sequences subject to recombination using parsimony. *Math Biosci*. 98:185–200. [https://doi.org/10.1016/0025-5564\(90\)90123-G](https://doi.org/10.1016/0025-5564(90)90123-G).
- Huang Z, Kelleher J, Chan Y-b, Balding DJ. 2024. Estimating evolutionary and demographic parameters via ARG-derived IBD. *PLoS Genet*. 21:e1011537.
- Hudson RR. 1983. Properties of a neutral allele model with intragenic recombination. *Theor Popul Biol*. 23:183–201. [https://doi.org/10.1016/0040-5809\(83\)90013-8](https://doi.org/10.1016/0040-5809(83)90013-8).
- Hudson RR, Kaplan NL. 1985. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*. 111:147–164. <https://doi.org/10.1093/genetics/111.1.147>.
- Ignatieva A, Favero M, Koskela J, Sant J, Myers SR. 2025. The length of haplotype blocks and signals of structural variation in reconstructed genealogies. *Mol Biol Evol*. 42:msaf190.
- Kelleher J et al. 2019. Inferring whole-genome histories in large population datasets. *Nat Genet*. 51:1330–1338. <https://doi.org/10.1038/s41588-019-0483-y>.
- Kelleher J et al. 2024. tskit-dev/tskit: Python 0.5.8.
- Kelleher J, Etheridge AM, McVean G. 2016. Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Comput Biol*. 12:e1004842. <https://doi.org/10.1371/journal.pcbi.1004842>.
- Kelleher J, Thornton KR, Ashander J, Ralph PL. 2018. Efficient pedigree recording for fast population genetics simulation. *PLoS Comput Biol*. 14:1–21. <https://doi.org/10.1371/journal.pcbi.1006581>.
- Kendall M, Colijn C. 2016. Mapping phylogenetic trees to reveal distinct patterns of evolution. *Mol Biol Evol*. 33:2735–2743. <https://doi.org/10.1093/molbev/msw124>.
- Kuhner MK, Felsenstein J. 1994. A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol Biol Evol*. 11:459–468. <https://doi.org/10.1093/oxfordjournals.molbev.a040126>.
- Kuhner MK, Yamato J. 2015. Assessing differences between ancestral recombination graphs. *J Mol Evol*. 80:258–264. <https://doi.org/10.1007/s00239-015-9676-x>.
- Kuhner MK, Yamato J, Felsenstein J. 2000. Maximum likelihood estimation of recombination rates from population data. *Genetics*. 156:1393–1401. <https://doi.org/10.1093/genetics/156.3.1393>.
- 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: 1–24. <https://doi.org/10.1371/journal.pgen.1011110>.
- Llabrés M, Rosselló F, Valiente G. 2021. The generalized robinson-foulds distance for phylogenetic trees. *J Comput Biol*. 28:1181–1195. PMID: 34714118. <https://doi.org/10.1089/cmb.2021.0342>.
- Minichiello MJ, Durbin R. 2006. Mapping trait loci by use of inferred ancestral recombination graphs. *Am J Hum Genet*. 79:910–922. <https://doi.org/10.1086/508901>.
- Myers SR, Griffiths RC. 2003. Bounds on the minimum number of recombination events in a sample history. *Genetics*. 163:375–394. <https://doi.org/10.1093/genetics/163.1.375>.
- Nielsen R, Vaughn AH, Deng Y. 2025. Inference and applications of ancestral recombination graphs. *Nat Rev Genet* 47–58. <https://doi.org/10.1038/s41576-023-00677-8>.
- Ralph P, Coop G. 2013. The geography of recent genetic ancestry across Europe. *PLoS Biol*. 11:e1001555. <https://doi.org/10.1371/journal.pbio.1001555>.
- Ralph P, Thornton K, Kelleher J. 2020. Efficiently summarizing relationships in large samples: a general duality between statistics of genealogies and genomes. *Genetics*. 215:779–797. <https://doi.org/10.1534/genetics.120.303253>.

- Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. 2014. Genome-wide inference of ancestral recombination graphs. *PLoS Genet.* 10: e1004342. <https://doi.org/10.1371/journal.pgen.1004342>.
- Ringbauer H, Coop G, Barton NH. 2017. Inferring recent demography from isolation by distance of long shared sequence blocks. *Genetics.* 205: 1335–1351. <https://doi.org/10.1534/genetics.116.196220>.
- Robinson DF, Foulds LR. 1979. Comparison of weighted labelled trees. In Horadam AF, Wallis WD, editors. *Combinatorial mathematics VI*. Berlin, Heidelberg. Springer Berlin Heidelberg. p. 119–126.
- Robinson D, Foulds L. 1981. Comparison of phylogenetic trees. *Math Biosci.* 53:131–147. [https://doi.org/10.1016/0025-5564\(81\)90043-2](https://doi.org/10.1016/0025-5564(81)90043-2).
- Silcocks M et al. 2023. Indigenous Australian genomes show deep structure and rich novel variation. *Nature.* 624:593–601. <https://doi.org/10.1038/s41586-023-06831-w>.
- Song YS, Hein J. 2005. Constructing minimal ancestral recombination graphs. *J Comput Biol.* 12:147–169. <https://doi.org/10.1089/cmb.2005.12.147>.
- Speidel L, Forest M, Shi S, Myers SR. 2019. A method for genome-wide genealogy estimation for thousands of samples. *Nat Genet.* 51: 1321–1329. <https://doi.org/10.1038/s41588-019-0484-x>.
- Stephens M, Donnelly P. 2000. Inference in molecular population genetics. *J R Stat Soc Series B Stat Methodol.* 62:605–635. <https://doi.org/10.1111/1467-9868.00254>.
- Tsambos G, Kelleher J, Ralph P, Leslie S, Vukcevic D. 2023. link-ancestors: fast simulation of local ancestry with tree sequence software. *Bioinform Adv.* 3:vbv163. <https://doi.org/10.1093/bioadv/vbv163>.
- Wang L, Zhang K, Zhang L. 2001. Perfect phylogenetic networks with recombination. *J Comput Biol.* 8:69–78. <https://doi.org/10.1089/106652701300099119>.
- Wohns AW et al. 2022. A unified genealogy of modern and ancient genomes. *Science.* 375:eabi8264. <https://doi.org/10.1126/science.abi8264>.
- Wong Y et al. 2024. A general and efficient representation of ancestral recombination graphs. *Genetics.* 228:iyae100. <https://doi.org/10.1093/genetics/iyae100>.
- Yang S, Carmi S, Pe'er I. 2016. Rapidly registering identity-by-descent across ancestral recombination graphs. *J Comput Biol.* 23: 495–507. <https://doi.org/10.1089/cmb.2016.0016>.
- Zhang BC, Biddanda A, Gunnarsson Á.F, Cooper F, Palamara PF. 2023. Biobank-scale inference of ancestral recombination graphs enables genealogical analysis of complex traits. *Nat Genet.* 55:768–776. <https://doi.org/10.1038/s41588-023-01379-x>.

Editor: S. Edwards