



Genetics and population analysis

KwARG: parsimonious reconstruction of ancestral recombination graphs with recurrent mutation

Anastasia Ignatieva ^{1,*}, Rune B. Lyngsø², Paul A. Jenkins ^{1,3,4} and Jotun Hein^{2,4}

¹Department of Statistics, University of Warwick, Coventry CV4 7AL, UK, ²Department of Statistics, University of Oxford, Oxford OX1 3LB, UK, ³Department of Computer Science, University of Warwick, Coventry CV4 7AL, UK and ⁴The Alan Turing Institute, British Library, London NW1 2DB, UK

*To whom correspondence should be addressed.

Associate Editor: Russell Schwartz

Received on December 18, 2020; revised on April 30, 2021; editorial decision on May 5, 2021; accepted on May 7, 2021

Abstract

Motivation: The reconstruction of possible histories given a sample of genetic data in the presence of recombination and recurrent mutation is a challenging problem, but can provide key insights into the evolution of a population. We present KwARG, which implements a parsimony-based greedy heuristic algorithm for finding plausible genealogical histories (ancestral recombination graphs) that are minimal or near-minimal in the number of posited recombination and mutation events.

Results: Given an input dataset of aligned sequences, KwARG outputs a list of possible candidate solutions, each comprising a list of mutation and recombination events that could have generated the dataset; the relative proportion of recombinations and recurrent mutations in a solution can be controlled via specifying a set of ‘cost’ parameters. We demonstrate that the algorithm performs well when compared against existing methods.

Availability and implementation: The software is available at <https://github.com/a-ignatieva/kwarg>.

Contact: anastasia.ignatieva@warwick.ac.uk

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

For many species, the evolution of genetic variation within a population is driven by the processes of mutation and recombination in addition to genetic drift. A typical mutation affects the genome at a single position, and may or may not spread through subsequent generations by inheritance. Recombination, on the other hand, occurs when a new haplotype is created as a mixture of genetic material from two different sources, which can drive evolution at a much faster rate. The detection of recombination is an important problem that can provide crucial scientific insights, for instance in understanding the potential for rapid changes in pathogenic properties within viral populations (Simon-Loriere and Holmes, 2011).

Consider a population evolving through the replication, mutation and recombination of genetic material within individuals, emerging from a common origin and living through multiple generations until the present day. In general, the history of shared ancestry, mutation and recombination events are not observed, and must be inferred from a sample of genetic data obtained from the present-day population. Crossover recombination can occur anywhere along a sequence, and the breakpoint position is also unobserved. This article focuses on methods for reconstructing possible histories of such a sample, in the form of *ancestral recombination graphs* (ARGs)—

networks of evolution connecting the sampled individuals to shared ancestors in the past through coalescence, mutation and crossover recombination events; an example is illustrated in Figure 1. This is a very important but challenging problem, as many possible histories might have generated a given sample. Moreover, recombination can be undetectable unless mutations appear on specific branches of the genealogy (Hein *et al.*, 2004, Section 5.11), and recombination events can produce patterns in the data that are indistinguishable from the effects of *recurrent mutation* (McVean *et al.*, 2002); that is, two or more mutation events in a genealogical history that affect the same locus.

Parsimony is an approach focused on finding possible histories which minimize the number of recombinations and recurrent mutations. This does not necessarily describe the most biologically plausible version of events, but produces a useful lower bound on the complexity of the evolutionary pathway that might have generated the given dataset. Beyond specifying the types of events that are allowed, parsimony does not require assuming a particular generative model; the approach focuses on sequences of events that can generate the observed dataset, disregarding the timing and prior rate of these events.

Previous work on reconstructing histories using parsimony has tackled recombination and recurrent mutation separately. Algorithms for reconstructing minimal ARGs generally make the

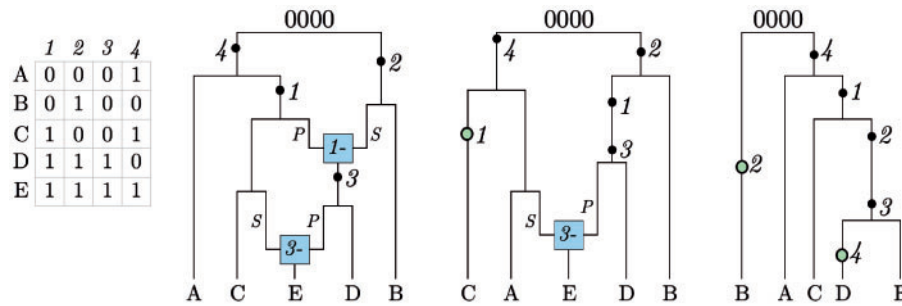


Fig. 1. Three examples of ARGs. The dataset is shown on the left in binary format, with 0's and 1's corresponding to the ancestral and mutant state at each site, respectively. Mutation events are shown as black dots and labelled by the site they affect; green filled circle corresponds to a recurrent mutation. Recombination nodes (in blue) are labelled with the recombination breakpoint; material to the right (left) of the breakpoint is inherited from the parent connected by the edge labelled S ('suffix') ('prefix')

infinite sites assumption, which allows at most one mutation to have occurred at each site of the genome, thus precluding recurrent mutation events, and the goal is to calculate the minimum number of crossover recombinations required to explain a dataset, denoted R_{\min} . Even with this constraint, the problem is NP-hard (Wang *et al.*, 2001); exact algorithms are practical only for small datasets (Hein, 1990; Lyngso *et al.*, 2005), and general methods rely on heuristic approximations (Hein, 1993; Minichiello and Durbin, 2006; Parida *et al.*, 2008; Song *et al.*, 2005; Thao and Vinh, 2019). Alternatively, one can assume the absence of recombination and seek to calculate the minimum number of recurrent mutations required, denoted P_{\min} . In this case, reconstruction of maximum parsimony trees is also NP-hard (Foulds and Graham, 1982); likewise, methods can only handle small datasets or are based on heuristics (Semple and Steel, 2003, Section 5.4).

Parsimony contrasts with the alternative approach of model-based inference, which requires the user to select a generative model and relies on the estimation of mutation and recombination rates as model parameters. Model-based inference generally involves integrating over the space of possible histories, which is usually intractable; methods rely on MCMC (e.g. Rasmussen *et al.*, 2014) or importance sampling (e.g. Jenkins and Griffiths, 2011), but the problem remains computationally difficult. If the presence of recombination is certain and reasonable models of population dynamics are available, model-based approaches may be more suitable and result in more powerful inference. However, model misspecification can play an important role, for instance when modelling viral evolution over a transmission network, where the relative importance of factors such as geographical structure, social clustering and the impact of interventions may be difficult to ascertain. In this case, model-based inference can provide misleading results if overinterpreted, with poor quantification of uncertainty due to model misspecification. Parsimony-based methods fail to offer the interpretability or uncertainty quantification of a model but this does preclude their results being overinterpreted. They are simple and straightforward to implement and can be useful in situations such as enabling testing for the presence or absence of recombination when this is not certain (Bruen *et al.*, 2006).

There are a number of recently developed methods, namely RENT+ (Mirzaei and Wu, 2017), tsinfer (Kelleher *et al.*, 2019) and Relate (Speidel *et al.*, 2019), that seek to reconstruct local tree or ARG topologies from the data. These methods do not make strict model-based assumptions, incorporating heuristic algorithms, and do not aim to reconstruct the most *parsimonious* histories. We note also the existence of numerous other methods for inference of recombination (e.g. Boni *et al.*, 2007; Kosakovsky Pond *et al.*, 2006; Li and Stephens, 2003; Martin and Rybicki, 2000) which do not explicitly reconstruct ARGs.

KwARG ('quick ARG') is a software tool, written in C, which implements a greedy heuristic-based parsimony algorithm for reconstructing histories that are minimal or near-minimal in the number of posited recombination and mutation events. The algorithm starts with the input dataset and generates plausible histories backwards in time, adding coalescence, mutation, recombination and recurrent mutation events to reduce the dataset until the common ancestor is reached. By tuning a set of cost parameters for each event type,

KwARG can find solutions consisting only of recombinations (giving an upper bound on R_{\min}), only of recurrent mutations (giving an upper bound on P_{\min}), or a combination of both event types. KwARG handles both the 'infinite sites' and 'maximum parsimony' scenarios, as well as interpolating between these two cases by allowing recombinations as well as recurrent mutations and sequencing errors, which is not offered by existing methods. This is illustrated in Figure 1: KwARG finds all three types of solution for the given dataset. KwARG shows excellent performance when benchmarked against exact methods on small datasets, and outperforms existing parsimony-based heuristic methods on large, more complex datasets while maintaining computational efficiency; KwARG also achieves very good accuracy in reconstructing local tree topologies. The source code and executables are made freely available on GitHub at <https://github.com/a-ignatieva/kwarg>, along with documentation and usage examples.

The article is structured as follows. Details of the algorithm underlying KwARG are given in Section 2, with an explanation of the required inputs and expected outputs. In Section 3, the performance of KwARG on simulated data is benchmarked against exact methods and existing programs. An application of KwARG to a widely studied *Drosophila melanogaster* dataset (Kreitman, 1983) is described in Section 4. Discussion follows in Section 5.

2 Materials and Methods

Consider a sample of genetic data, where the allele at each site can be denoted 0 or 1. We do not make the infinite sites assumption, so that each site can undergo multiple mutation events. However, we do assume that mutations correspond to transitions between exactly two possible states, excluding for instance triallelic sites.

2.1 Input

KwARG accepts data in the form of a binary matrix, or a multiple alignment in nucleotide or amino acid format. The sequence and site labels can be provided if desired. It is possible to specify a root sequence, or leave this to be determined. The presence of missing data are permitted; regardless of the type of input, the data are converted to a binary matrix \mathcal{D} , with entries '*' denoting missing entries or material that is not ancestral to the sample.

2.2 Methods

Under the infinite sites assumption, at most one mutation is allowed to have occurred per site. If any two columns contain all four of the configurations 00, 01, 10, 11, then the data could not have been generated only through replication and mutation, and there must have been at least one recombination event between the two corresponding sites. This is the four gamete test (Hudson and Kaplan, 1985), and the two sites are said to be *incompatible*. When recurrent mutations are allowed, the incompatibility could likewise have been generated through multiple mutations affecting the same site (McVean *et al.*, 2002).

KwARG reconstructs the history of a sample backwards in time, by starting with the data matrix \mathcal{D} and performing row and column operations corresponding to coalescence, mutation and recombination events, until only one ancestral sequence remains. By reversing the order of the steps, a forward-in-time history is obtained, showing how the population evolved from the ancestor to the present sample. When a choice can be made between multiple possible events, a neighbourhood of candidate ancestral states is constructed, using the same general method as that employed in the program Beagle (Lyngsø *et al.*, 2005). A backwards-in-time approach has also been implemented in the programs SHRUB (Song *et al.*, 2005), Margarita (Minichiello and Durbin, 2006) and GAMARG (Thao and Vinh, 2019), all of which adopt the infinite sites assumption but use different criteria for choosing amongst possible recombination events.

2.2.1 Construction of a history

For convenience, assume that the all-zero sequence is specified as the root, and 0 (1) entries of \mathcal{D} correspond to ancestral (mutated) sites. Suppose \mathcal{D}_t is the data matrix obtained after $t - 1$ iterations of the algorithm. At the beginning of the t th step, KwARG first reduces \mathcal{D}_t , by repeatedly applying the ‘Clean’ algorithm (Song and Hein, 2003) through:

- deleting uninformative columns (consisting of all 0’s);
- deleting columns containing only one 1 (corresponding to ‘undoing’ a mutation present in only one sequence);
- deleting a row if it agrees with another row (corresponding to a coalescence event);
- deleting a column if it agrees with an adjacent column.

Two rows (columns) *agree* if they are equal at all positions where both rows (columns) contain ancestral material, and the sites (sequences) carrying ancestral material in one are a subset of the sites (sequences) carrying ancestral material in the other.

A run of the ‘Clean’ algorithm repeatedly applies these steps to \mathcal{D}_t , terminating when no further reduction is possible. Suppose the resulting data matrix is $\overline{\mathcal{D}}_t$. KwARG then constructs a neighbourhood \mathcal{N}_t of candidate next states, each one obtained through one of the following operations:

- Pick a row and split it into two at a possible recombination point. Only a subset of possible recombining sequences and breakpoints needs to be considered; see Lyngsø *et al.* (2005, Section 3.3) for a detailed explanation.
- Remove a recurrent mutation, by selecting a column and changing a 0 entry to 1, or a 1 entry to 0. This is the event type that is disallowed by algorithms applying the infinite sites assumption.

Suppose a neighbourhood $\mathcal{N}_t = \{\mathcal{N}_t^1, \dots, \mathcal{N}_t^N\}$ is formed, consisting of all possible states that can be reached from $\overline{\mathcal{D}}_t$ through applying one of these operations. Then the reduced neighbourhood $\overline{\mathcal{N}}_t = \{\overline{\mathcal{N}}_t^1, \dots, \overline{\mathcal{N}}_t^N\}$ is formed by applying ‘Clean’ to each state in turn. Each state $\overline{\mathcal{N}}_t^i$ is then assigned a score $S(\overline{\mathcal{N}}_t^i, \mathcal{N}_t^i, \overline{\mathcal{D}}_t)$, combining (i) the cost $C(\mathcal{N}_t^i, \overline{\mathcal{D}}_t)$, defined below, of reaching the configuration \mathcal{N}_t^i from $\overline{\mathcal{D}}_t$, (ii) a measure $\text{AM}(\overline{\mathcal{N}}_t^i)$ of the complexity of the resulting data matrix $\overline{\mathcal{N}}_t^i$ and (iii) a lower bound $L(\overline{\mathcal{N}}_t^i)$ on the remaining number of recombination and recurrent mutation events still required to reach the ancestral sequence from $\overline{\mathcal{N}}_t^i$. Finally, a state is selected, say $\overline{\mathcal{N}}_t^i$, based on its score, and we set $\mathcal{D}_{t+1} = \overline{\mathcal{N}}_t^i$. The process of reducing the dataset followed by constructing a neighbourhood and choosing the best move is repeated, until all incompatibilities are resolved and the root sequence is reached. Pseudocode for the ‘Clean’ algorithm and KwARG is given in [Supplementary Section S1](#).

The construction of a history for the dataset given in [Figure 1](#) is illustrated in [Figure 2](#). The first step corresponds to the construction of a neighbourhood, two of the states $\mathcal{N}_1^1, \mathcal{N}_1^2 \in \mathcal{N}_1$ are pictured.

Then, the ‘Clean’ algorithm is applied to each state in the neighbourhood (illustrated as a series of steps following blue arrows). From the resulting reduced neighbourhood $\{\overline{\mathcal{N}}_1^1, \overline{\mathcal{N}}_1^2, \dots\}$, the state $\overline{\mathcal{N}}_1^2$ is selected; the other illustrated path is abandoned. This process is repeated until all incompatibilities are resolved and the empty state is reached. Following the path of selected moves in this figure left-to-right corresponds to the events encountered when traversing the leftmost ARG in [Figure 1](#) from the bottom up. If instead the state $\overline{\mathcal{N}}_1^1$ were selected at the second step of the algorithm, the resulting path would correspond to the ARG in the centre of [Figure 1](#).

2.2.2 Score

When considering which next step to take, more informed choices can be made by considering not just the cost of the step, but also the complexity of the configuration it leads to. This is the principle behind the A* algorithm (Hart *et al.*, 1968), using a heuristic estimate of remaining distance to guide the choice of the next node to expand. KwARG applies the same principle in a greedy fashion, following a path of locally optimal choices in an attempt to find a minimal history.

The score implemented in KwARG is

$$S(\overline{\mathcal{N}}_t^i, \mathcal{N}_t^i, \overline{\mathcal{D}}_t) = \left(C(\mathcal{N}_t^i, \overline{\mathcal{D}}_t) + L(\overline{\mathcal{N}}_t^i) \right) \cdot \max\text{AM}(\overline{\mathcal{N}}_t^i) + \text{AM}(\overline{\mathcal{N}}_t^i), \quad (1)$$

where

$$L(\overline{\mathcal{N}}_t^i) = \begin{cases} R_{\min}(\overline{\mathcal{N}}_t^i) & \text{if } \max\text{AM}(\overline{\mathcal{N}}_t^i) < 75, \\ HB(\overline{\mathcal{N}}_t^i) & \text{if } 75 \leq \max\text{AM}(\overline{\mathcal{N}}_t^i) < 200, \\ HK(\overline{\mathcal{N}}_t^i) & \text{otherwise.} \end{cases}$$

Here, $C(\mathcal{N}_t^i, \overline{\mathcal{D}}_t)$ denotes the cost of the corresponding event, defined in Section 2.2.3; $\max\text{AM}(\overline{\mathcal{N}}_t^i)$ denotes the maximum amount of ancestral material seen in any of the states in $\overline{\mathcal{N}}_t^i$, and $\text{AM}(\overline{\mathcal{N}}_t^i)$ gives the amount of ancestral material in state $\overline{\mathcal{N}}_t^i$. Incorporating a measure of the amount of ancestral material in a state helps to break ties by assigning a smaller score to simpler configurations.

The method of computing the lower bound L depends on the complexity of the dataset, with a trade-off between accuracy and computational cost. For relatively small datasets, it is feasible to compute R_{\min} exactly using Beagle. HB refers to the haplotype bound, employing the improvements afforded by first calculating local bounds for incompatible intervals, and applying a composition method to obtain a global bound (Myers and Griffiths, 2003). HK refers to the Hudson–Kaplan bound (Hudson and Kaplan, 1985); this is quick but less accurate, so is reserved for larger, more complex configurations. Note that these bounds are computed under the infinite sites assumption.

The particular form and components of the score were chosen through simulation testing; we found that the given formula provides a good level of informativeness regarding the quality of a possible state.

2.2.3 Event cost

Each type of event is assigned a cost, which gives a relative measure of preference for each event type in the reconstructed history:

- C_R : the cost of a single recombination event, defaults to 1.
- C_{RR} : the cost of performing two successive recombinations, defaults to 2. It is sufficient to consider at most two consecutive recombination events before a coalescence (Lyngsø *et al.*, 2005); this type of event also captures the effects of gene conversion.
- C_{RM} : the cost of a recurrent mutation. If \mathcal{N}_t^i is formed from $\overline{\mathcal{D}}_t$ by a recurrent mutation in a column representing k agreeing sites, this corresponds to proposing k recurrent mutation events, so the cost is $C(\mathcal{N}_t^i, \overline{\mathcal{D}}_t) = k \cdot C_{RM}$.

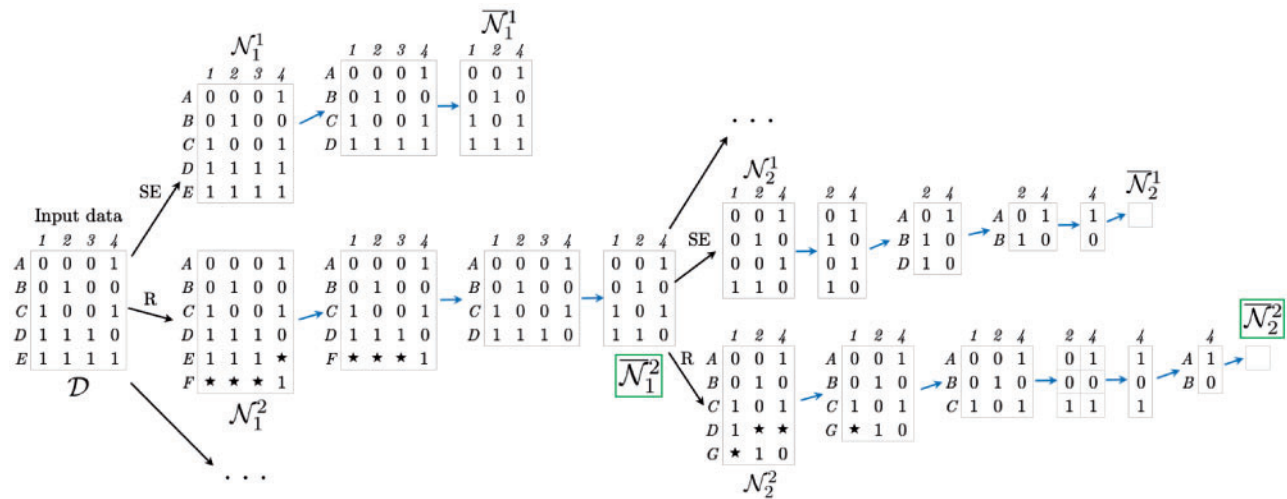


Fig. 2. Example of a reconstructed history for the dataset in Figure 1. Stars ‘*’ denote non-ancestral material. SE: recurrent mutation occurring on a terminal branch of the ARG. R: recombination event. A sequence of blue arrows corresponds to one application of the ‘Clean’ algorithm. Green boxes highlight the selected states

- C_{SE} : this event is a recurrent mutation which affects only one sequence in the original dataset, i.e. it occurs on the terminal branches of the ARG. Thus, the event can be either a regular recurrent mutation or an artefact due to sequencing errors. The cost can be set to equal C_{RM} , or lower if the presence of sequencing errors is considered likely.

KwARG allows the specification of a range of event costs as tuning parameters, as well as the number Q of independent runs of the algorithm to perform for each cost configuration. The proportions of recombinations to recurrent mutations in the solutions produced by KwARG can be controlled by varying the ratio of costs for the corresponding event types.

2.2.4 Selection probability

The method of selecting the next state from a neighbourhood of candidates will impact on the efficiency and performance of the algorithm. At one extreme, selecting at random amongst the states will mean that the solution space is explored more fully, but will be prohibitively inefficient in terms of the number of runs needed to find a near-optimal solution. On the other hand, always greedily selecting the move with the minimal score will quickly identify a small set of solutions for each cost configuration, at the expense of placing our faith in the ability of the score to assess the quality of the candidate states accurately.

We propose a selection method that is intermediate between these two extremes, randomizing the selection but focusing on moves with near-minimal scores. A pseudo-score for state \bar{N}_t^i is calculated:

$$\exp\left(T \cdot \left(1 - \tilde{S}(\bar{N}_t^i, N_t^i, \bar{D}_t)\right)\right), \quad (2)$$

where

$$\tilde{S}(\bar{N}_t^i, N_t^i, \bar{D}_t) = \frac{S(\bar{N}_t^i, N_t^i, \bar{D}_t) - \min_j S(\bar{N}_t^j, N_t^j, \bar{D}_t)}{\max_j S(\bar{N}_t^j, N_t^j, \bar{D}_t) - \min_j S(\bar{N}_t^j, N_t^j, \bar{D}_t)},$$

and states in \bar{N}_t are selected with probability proportional to their pseudo-score. The annealing parameter T controls the extent of random exploration; $T=0$ corresponds to choosing uniformly at random from the neighbourhood of candidates, and $T=\infty$ to always choosing a state with the minimal score. The default value of $T=30$ was chosen following simulation testing, which showed that this

provides a good balance between efficiency and thorough exploration of the neighbourhood.

2.3 Output

The default output consists of the number of recombinations and recurrent mutations in each identified solution; an example for the Kreitman dataset is given in Table 1. Each iteration is assigned a unique random seed, which can be used to reconstruct each particular solution and produce more detailed outputs, such as a detailed list of events in the history, the ARG in several graph formats or the corresponding sequence of marginal trees.

3 Results

We have tested the performance of KwARG on simulated data, based on two main criteria. First, we compared its performance against exact methods, PAUP* and Beagle, to demonstrate that KwARG successfully reconstructs minimal histories in the mutation-only and recombination-only cases, respectively. Second, we carried out simulation studies to determine how accurately KwARG reconstructs local trees, compared against three other methods: tsinfer, RENT+ and ARGweaver. Finally, we compared how well KwARG performs against the parsimony-based heuristic methods SHRUB (Song et al., 2005) and SHRUB-GC (Song et al., 2006); these results are presented in Supplementary Section S4. We also investigated the dependence of the run time of KwARG on the number and length of sequences, through simulation studies.

Table 1. Example output of KwARG for the Kreitman dataset

Seed	T	C_{SE}	C_{RM}	C_R	C_{RR}	SE	RM	R	$\sum_t \bar{N}_t^i $
2263536315	30.0	∞	∞	1.00	2.00	0	0	7	143
2347021759	30.0	0.90	0.91	1.00	2.00	1	0	6	853
1791455164	30.0	0.80	0.81	1.00	2.00	1	0	5	728
1684879495	30.0	0.60	0.61	1.00	2.00	2	0	4	783
1884182000	30.0	0.40	0.41	1.00	2.00	3	0	3	806
1900122424	30.0	0.20	0.21	1.00	2.00	5	0	2	702
2111915557	30.0	0.10	0.11	1.00	2.00	8	0	1	833
2888657821	30.0	0.01	0.02	1.00	2.00	10	0	0	715

Note: SE: number of recurrent mutations occurring on terminal branches of the ARG (possible sequencing errors); RM: number of other recurrent mutations; R: number of recombinations. Last column gives the total number of neighbourhood states considered.

3.1 Finite sites

3.1.1 Comparison to PAUP*

Disallowing recombination, the quality of computed upper bounds on P_{\min} was tested by comparison with PAUP* (Swofford, 2001, version 4.0a168), which was used to compute the exact minimum parsimony score via branch-and-bound on 994 datasets simulated as described in Supplementary Section S3.1.

KwARG failed to find P_{\min} in 11 (1.1%) cases out of 994. The results are illustrated in the top panel of Figure 3. Where KwARG failed to find an optimal solution, in all 11 cases it was off by just one recurrent mutation. Figure 3 also demonstrates that a substantial proportion of recurrent mutations do not create incompatibilities in the data, and the number of actual events often far exceeds P_{\min} .

3.2 Infinite sites

3.2.1 Comparison to beagle

Under the infinite sites assumption (disallowing recurrent mutation), the accuracy of KwARG's upper bound on R_{\min} was tested by comparison with Beagle (Lyngsø et al., 2005), on 1037 datasets simulated as described in Supplementary Section S3.2.

Using the default annealing parameter $T = 30$, KwARG found R_{\min} in all cases. In 97% of the runs, this took under 5 s of CPU time (on a 2.7 GHz Intel Core i7 processor); all but one run took

<40 s. In 93% of the runs, one iteration was sufficient to find an optimal solution; in 99% of the runs, five iterations were sufficient. Beagle found the exact solution in 5 s or less in 86% of cases; for datasets with a small R_{\min} Beagle runs relatively quickly (median run time for $R_{\min} = 5$ was 1 s, compared to KwARG's 0.3 s). For more complex datasets, KwARG finds an optimal solution much faster; for $R_{\min} = 9$, the median run time of Beagle was 56 s, compared to KwARG's 3 s.

Setting $T = 10$ and $T = \infty$ resulted in 5 and 22 failures to find an optimal solution, respectively, when KwARG was run for $Q = 1000$ iterations per dataset (or terminated after 10 min have elapsed), demonstrating that setting the annealing parameters too low or too high results in deterioration of performance.

The bottom panel of Figure 3 illustrates the results and shows the relationship between the true simulated number of recombinations and R_{\min} . This demonstrates that in many cases, substantially more recombinations have occurred than can be confidently detected from the data.

3.2.2 Comparison to tsinfer, RENT+ and ARGweaver

We tested the performance of KwARG in recovering the topology of simulated local trees for a range of recombination and mutation rates (under the infinite sites assumption). For each combination of rates, we simulated 100 datasets; details of the simulation parameters and settings used in running each program are given in Supplementary Section S5. From the output of each method, we calculated the Kendall–Colijn metric (Kendall and Colijn, 2016) between the inferred and true tree topologies at each variant site position, calculating the mean across all variant sites and averaging over the 100 datasets. We note that ARGs contain more information than local trees, but there is no obvious way of comparing ARG topologies (and tsinfer only infers local trees, rather than full ARGs).

The results are shown in the top panel of Figure 4 and Supplementary Figure S4. All methods show very comparable performance across the range of considered scenarios, with KwARG slightly outperforming the other methods, based on the chosen metric, when the recombination rate is relatively low and the mutation rate relatively high. We have performed the same analysis using the Robinson–Foulds metric (Robinson and Foulds, 1981), and found this to give very similar results.

3.3 Run time analysis

A comparison of the run times of KwARG against tsinfer, RENT+ and ARGweaver is presented in the bottom panel of Figure 4 and Supplementary Figure S5. KwARG demonstrates good efficiency when the recombination and mutation rates are relatively low, and shows roughly linear growth in run time as the mutation rate increases.

The dependence of the run time of KwARG on the number and length of sequences was further investigated through simulations; the results are presented in Supplementary Section S6. Keeping the sequence length fixed showed that KwARG runs very quickly when the number of sequences is very low, and shows roughly exponential growth in run time when the number of sequences is 6 or more. Keeping the number of sequences fixed shows that, after an initial exponential increase (due to small datasets taking very little time per iteration), the run time scales roughly linearly in sequence length.

3.4 Application to Kreitman data

The performance of KwARG is illustrated on the classic dataset of Kreitman (1983, Table 1); this is not close to the performance limit of KwARG, but has been widely used for benchmarking algorithms used for ARG reconstruction. The dataset consists of 11 sequences and 2721 sites, of which 43 are polymorphic, of the alcohol dehydrogenase locus of *D.melanogaster*. The data are shown in Figure 5, with columns containing singleton mutations removed for ease of viewing. Applying the 'Clean' algorithm, as described in Section 2.2.1, reduces this to matrix of 9 rows and 16 columns. KwARG was run with the default parameters, $Q = 500$ times for each of 13

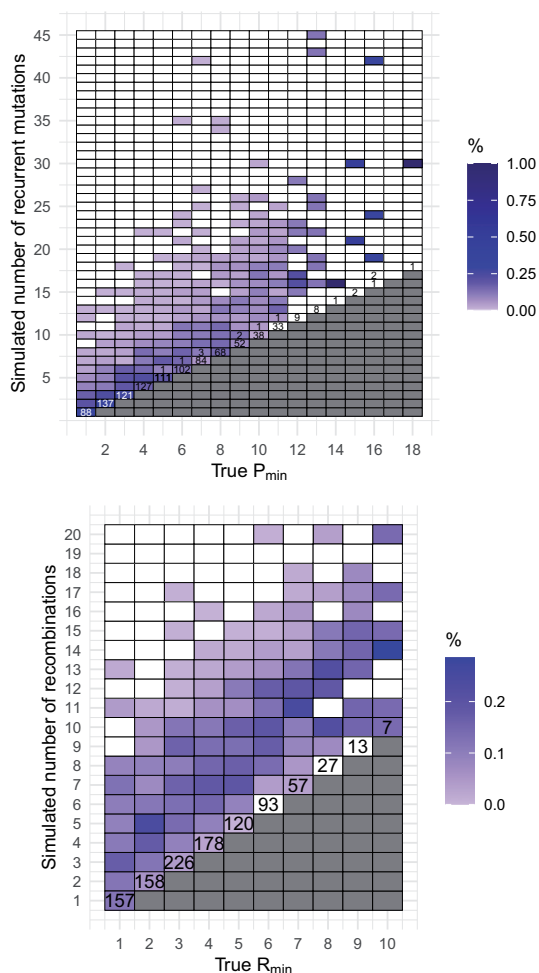


Fig. 3. Top: number of simulated recurrent mutations against P_{\min} . Bottom: number of simulated recombinations against R_{\min} . Cell colouring intensity is proportional to the number of datasets generated for each pair of coordinates. Numbers in each cell correspond to the number of cases where for a dataset with the true minimum number of events given on the x-axis, KwARG inferred the number of events given on the y-axis (unlabelled cells correspond to 0 such cases)

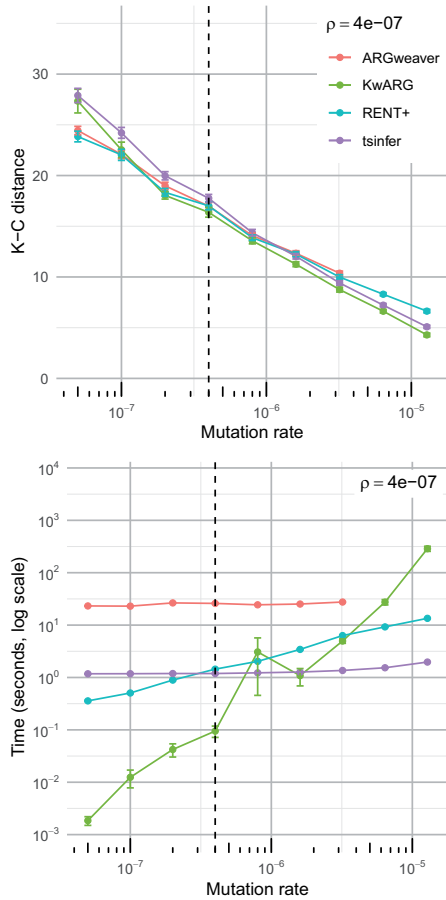


Fig. 4. Comparison of performance in inferring local trees. Top panel: points show mean across 100 simulated datasets for each value of mutation rate μ (per generation per site) with recombination rate $\rho = 4 \cdot 10^{-7}$ (per generation per site); error bars show mean \pm standard error. Lower K-C distance indicates better accuracy. Bottom panel: points show mean run time averaged over 100 datasets for each combination of rate parameters; error bars show mean \pm standard error. ARGweaver results not shown past $\mu = 3.2 \cdot 10^{-6}$ due to prohibitively long run time

default cost configurations given in Supplementary Section S2. An example of the output is shown in Table 1.

KwARG correctly identified the R_{\min} of 7 and the P_{\min} of 10 (confirmed by running Beagle and PAUP*, respectively). The 6500 iterations of KwARG took just under 9 min to run. Of these, 1829 (28%) resulted in optimal solutions; some are shown in Table 1.

KwARG identified multiple combinations of recombinations and recurrent mutations that could have generated this dataset. By default, slightly cheaper costs are assigned to recurrent mutations if they happen on terminal branches, so the results show a bias towards solutions with more SE events for each given number of recombinations.

The ten recurrent mutations appearing in the solution in row 8 of Table 1 are highlighted on the dataset in Figure 5. It is striking that 7 of these 10 recurrent mutations affect the same sequence FL-2S. In fact, these seven recurrent mutations could be replaced by three recombination events affecting sequence FL-2S, with breakpoints just after sites 3, 16 and 35; leaving the other identified recurrent mutations unchanged yields the solution in row 5 of Table 1. These findings suggest that the sequence may have been affected by cross-contamination or other errors during the sequencing process, or it could indeed be a recombinant mosaic of four other sequences in the sample. This recovers the results obtained by Stephens and Nei (1985), who posited the recombinant origins of sequence FL-2S following manual examination of a reconstructed maximum parsimony tree, which also highlighted the five consecutive mutations identified by KwARG. The ARG corresponding to the solution in row 5 of Table 1, visualized using Graphviz (Ellson et al., 2004), is shown in Figure 6.

Examination of the identified solutions also shows that site 36 of sequence Ja-S ‘necessitates’ two of the seven recombinations inferred in the minimal solution in the absence of recurrent mutation, while sites three and nine in sequences Wa-S and FL-1S, respectively, each create incompatibilities that could be resolved by one recombination.

4 Discussion

Methods for the reconstruction of parsimonious ARGs generally rely on the infinite sites assumption. When examining the output ARGs, it is often difficult to tell by how much the inferred recombination events actually affect the recombining sequences. As is the case with the Kreitman dataset, sometimes further examination reveals that two crossover recombination events have the same effect as one recurrent mutation, raising questions about which version of events is more likely. KwARG removes the need for such manual examination, and provides an automated way of highlighting such cases, which is particularly useful for larger datasets.

While KwARG performs well in inferring ARGs under the infinite sites assumption, it can be particularly useful in analysing genetic data from organisms whose genomes are reasonably likely to undergo recurrent mutation, such as viruses with relatively high mutation rates and short genomes. One such application is demonstrated in Ignatieva et al. (2021), where the output of KwARG is combined with probabilistic arguments to investigate the presence of ongoing recombination in SARS-CoV-2.

Zeros correspond to:

	C	C	C	A	A	G	G	C	G	A	C	C	C	C	G	G	A	T	C	T	C	T	A	T	T	C	G	C	C	
Wa-S	0	0	1	0	0	0	0	0	1	0	0	1	1	1	1	1	1	1	1	1	0	1	0	1	0	1	0	1	0	
FL-1S	0	0	0	0	0	1	0	0	0	1	0	0	1	1	1	1	1	1	1	1	0	1	0	1	0	1	0	1	0	
Af-S	0	0	1	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1	1	0
Fr-S	0	0	1	0	0	1	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1	1	0
FL-2S	0	0	1	1	1	0	1	1	0	0	1	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0
Ja-S	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
FL-F	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0
Fr-F	1	1	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	1
Wa-F	1	1	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	1
Af-F	1	1	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	1
Ja-F	1	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0

Fig. 5. Illustration of the Kreitman dataset. The 11 sequences labelled as in Kreitman (1983); polymorphic sites are labelled 1–43 and columns with singleton mutations are not shown

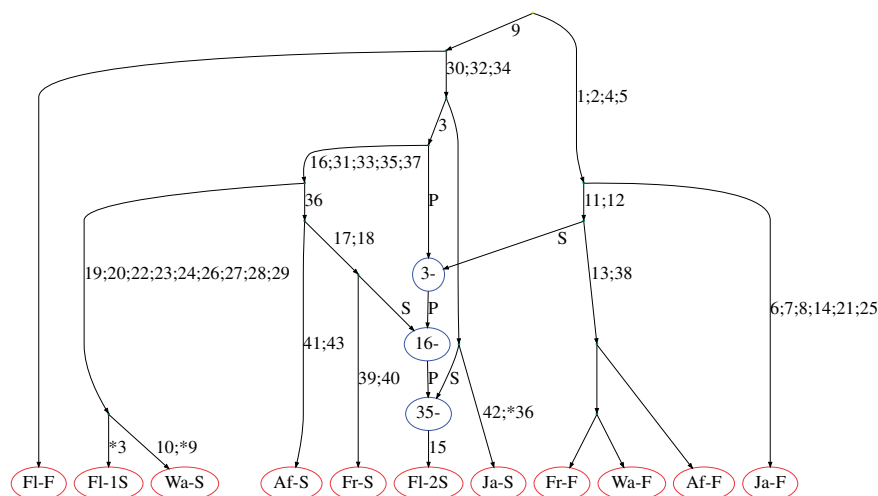


Fig. 6. ARG constructed for the Kreitman data. Edges are labelled with sites undergoing mutations; recurrent mutations are prefixed with an asterisk. Recombination nodes, in blue, are labelled with the recombination breakpoint; material to the right (left) of the breakpoint is inherited from the parent connected by the edge labelled *S* (for 'suffix') ('prefix')

The solutions identified by KwARG differ in the proportion of recurrent mutations to recombinations, ranging from an explanation that invokes only recombination events to one that invokes only mutation events. As is the case with other heuristic and parsimony-based methods, KwARG cannot offer uncertainty quantification for the inferred ARGs. Quantifying the likelihood of each scenario will be application-specific; for instance, one can choose a reasonable model of evolution for the population being studied, and identify the most likely solution under a range of reasonable mutation and recombination rates. When the presence or absence of recombination is not certain, then should the number of recurrent mutations needed to explain the dataset be infeasibly large, this provides evidence for the presence of recombination; this is the idea underlying the homoplasy test of Maynard Smith and Smith (1998). If the largest 'reasonable' number of recurrent mutations is then estimated, KwARG can be used to say how many additional recombination events are required to explain the dataset.

KwARG performs well when compared against exact parsimony methods for the 'recombination-only' and 'mutation-only' scenarios. Because of the random exploration incorporated within KwARG, it should be run multiple times on the same dataset before selecting the best solutions; the optimal run length of KwARG will be constrained by timing and the available computational resources. To gauge whether KwARG has run enough iterations, one could proceed by calculating R_{\min} and P_{\min} either exactly (if the data is reasonably small) or using other heuristics-based methods (such as SHRUB or PAUP*), to confirm whether KwARG has found good solutions at these two extremes.

The range of solutions explored by KwARG is guided by the choice of cost parameters. As a rule of thumb, simulations have shown that if the mutation and recombination rates are similar, costs near one give good accuracy of solutions in terms of reconstructing local tree topologies; if the mutation rate is significantly higher (lower) than the recombination rate, the cost should be set to less than (greater than) one. As KwARG incorporates a degree of random exploration, a range of solutions will still be obtained; the best choice of parameters will depend strongly on the nature and aims of the analysis being performed.

For model-based inference, the modelling assumptions can obviously affect the quality of the results; however, a parsimony-based approach also makes the strong assumption that the minimal ARG can capture useful information about the history of a sample. This will obviously depend strongly on the true recombination rate. Based on our comparisons with RENT+, tsinfer and ARGweaver, KwARG achieves very good accuracy of inference of local tree topologies at least comparable to these other methods, particularly when the recombination rate is low to moderate and the mutation

rate moderate to high. We emphasize that KwARG demonstrates relatively good accuracy even when the recombination rate is high and even though its express goal is to seek the most parsimonious, rather than necessarily the most likely, history. Moreover, for datasets with relatively few incompatibilities, the run time of KwARG is competitive with that of the other methods. It is also interesting to note that although all four programs incorporate very different approaches and heuristic algorithms, they demonstrate very similar performance in inferring local tree topologies over the range of considered scenarios.

The scalability of KwARG remains a challenge for large and more complex datasets. Performance gains could be readily achieved by running multiple iterations of KwARG in parallel, or incorporating more efficient ways of storing the intermediate states. Further improvements could also be obtained by amending the calculation of lower bounds within the cost function in order to account for the presence of recurrent mutation, which should make the scores more accurate, and hence the neighbourhood exploration more efficient. Other avenues for further work include explicitly incorporating gene conversion as a possible type of recombination event with a separate cost parameter, with a view to developing the underlying model of evolution to even more closely reflect biological reality.

Acknowledgements

We thank two anonymous reviewers for their helpful comments.

Funding

This work was supported by the Engineering and Physical Sciences Research Council and the Medical Research Council through the OxWaSP Centre for Doctoral Training [EPSRC grant number EP/L016710/1] and by the Alan Turing Institute [EPSRC grant number EP/N510129/1].

Conflict of Interest: none declared.

References

- Boni, M.F. *et al.* (2007) An exact nonparametric method for inferring mosaic structure in sequence triplets. *Genetics*, **176**, 1035–1047.
- Bruen, T.C. *et al.* (2006) A simple and robust statistical test for detecting the presence of recombination. *Genetics*, **172**, 2665–2681.
- Ellson, J. *et al.* (2004). Graphviz and Dynagraph: static and dynamic graph drawing tools. In: *Graph Drawing Software*, Springer, pp. 127–148.
- Foulds, L.R. and Graham, R.L. (1982) The Steiner problem in phylogeny is NP-complete. *Adv. Appl. Math.*, **3**, 43–49.

- Hart, P.E. et al. (1968) A formal basis for the heuristic determination of minimum cost paths. *IEEE Trans. Syst. Sci. Cybern.*, **4**, 100–107.
- Hein, J. (1990) Reconstructing evolution of sequences subject to recombination using parsimony. *Math. Biosci.*, **98**, 185–200.
- Hein, J. (1993) A heuristic method to reconstruct the history of sequences subject to recombination. *J. Mol. Evol.*, **36**, 396–405.
- Hein, J. et al. (2004). *Gene Genealogies, Variation and Evolution: A Primer in Coalescent Theory*. Oxford University Press, Oxford.
- Hudson, R.R. and Kaplan, N.L. (1985) Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*, **111**, 147–164.
- Ignatieva, A. et al. (2021) Investigation of ongoing recombination through genealogical reconstruction for SARS-CoV-2. *bioRxiv*.
- Jenkins, P.A. and Griffiths, R.C. (2011) Inference from samples of DNA sequences using a two-locus model. *J. Comput. Biol.*, **18**, 109–127.
- Kelleher, J. et al. (2019) Inferring whole-genome histories in large population datasets. *Nat. Genet.*, **51**, 1330–1338.
- Kendall, M. and Colijn, C. (2016) Mapping phylogenetic trees to reveal distinct patterns of evolution. *Mol. Biol. Evol.*, **33**, 2735–2743.
- Kosakovsky Pond, S.L. et al. (2006) Automated phylogenetic detection of recombination using a genetic algorithm. *Mol. Biol. Evol.*, **23**, 1891–1901.
- Kreitman, M. (1983) Nucleotide polymorphism at the alcohol dehydrogenase locus of *Drosophila melanogaster*. *Nature*, **304**, 412–417.
- Li, N. and Stephens, M. (2003) Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics*, **165**, 2213–2233.
- Lyngsø, R.B. et al. (2005). Minimum recombination histories by branch and bound. In: *International Workshop on Algorithms in Bioinformatics*, Springer, pp. 239–250.
- Martin, D. and Rybicki, E. (2000) RDP: detection of recombination amongst aligned sequences. *Bioinformatics*, **16**, 562–563.
- Maynard Smith, J. and Smith, N.H. (1998) Detecting recombination from gene trees. *Mol. Biol. Evol.*, **15**, 590–599.
- McVean, G. et al. (2002) A coalescent-based method for detecting and estimating recombination from gene sequences. *Genetics*, **160**, 1231–1241.
- Minichiello, M.J. and Durbin, R. (2006) Mapping trait loci by use of inferred ancestral recombination graphs. *Am. J. Hum. Genet.*, **79**, 910–922.
- Mirzaei, S. and Wu, Y. (2017) RENT+: an improved method for inferring local genealogical trees from haplotypes with recombination. *Bioinformatics*, **33**, 1021–1030.
- Myers, S.R. and Griffiths, R.C. (2003) Bounds on the minimum number of recombination events in a sample history. *Genetics*, **163**, 375–394.
- Parida, L. et al.; Genographic Consortium (2008) Estimating the ancestral recombinations graph (ARG) as compatible networks of SNP patterns. *J. Comput. Biol.*, **15**, 1133–1153.
- Rasmussen, M.D. et al. (2014) Genome-wide inference of ancestral recombination graphs. *PLoS Genet.*, **10**, e1004342.
- Robinson, D.F. and Foulds, L.R. (1981) Comparison of phylogenetic trees. *Math. Biosci.*, **53**, 131–147.
- Semple, C. and Steel, M. (2003) *Phylogenetics*. Oxford Lecture Series in Mathematics and its Applications, vol. 24. Oxford University Press, Oxford.
- Simon-Loriere, E. and Holmes, E.C. (2011) Why do RNA viruses recombine? *Nat. Rev. Microbiol.*, **9**, 617–626.
- Song, Y.S. and Hein, J. (2003) Parsimonious reconstruction of sequence evolution and haplotype blocks. In: *International Workshop on Algorithms in Bioinformatics*, Springer, pp. 287–302.
- Song, Y.S. et al. (2005) Efficient computation of close lower and upper bounds on the minimum number of recombinations in biological sequence evolution. *Bioinformatics*, **21**, i413–i422.
- Song, Y.S. et al. (2006) Algorithms to distinguish the role of gene-conversion from single-crossover recombination in the derivation of SNP sequences in populations. In: *Annual International Conference on Research in Computational Molecular Biology*, Springer, pp. 231–245.
- Speidel, L. et al. (2019) A method for genome-wide genealogy estimation for thousands of samples. *Nat. Genet.*, **51**, 1321–1329.
- Stephens, J.C. and Nei, M. (1985) Phylogenetic analysis of polymorphic DNA sequences at the ADH locus in *Drosophila melanogaster* and its sibling species. *J. Mol. Evol.*, **22**, 289–300.
- Swofford, D.L. (2003) *PAUP*: Phylogenetic Analysis Using Parsimony (and Other Methods)*. Version 4. Sinauer Associates, Sunderland, Massachusetts.
- Thao, N.T.P. and Vinh, L.S. (2019) A hybrid approach to optimize the number of recombinations in ancestral recombination graphs. In: *Proceedings of the 2019 9th International Conference on Bioscience, Biochemistry and Bioinformatics*, pp. 36–42.
- Wang, L. et al. (2001) Perfect phylogenetic networks with recombination. *J. Comput. Biol.*, **8**, 69–78.