




The Number and Pattern of Viral Genomic Reassortments are not Necessarily Identifiable from Segment Trees

Qianying Lin ^{1,*} Emma E. Goldberg,¹ Thomas Leitner ¹ Carmen Molina-París ¹,
Aaron A. King ^{2,3,4,5} and Ethan O. Romero-Severson ^{1,*}

¹Theoretical Biology and Biophysics Group, Los Alamos National Laboratory, Los Alamos, NM, USA

²Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, MI, USA

³Department of Mathematics, University of Michigan, Ann Arbor, MI, USA

⁴Center for the Study of Complex Systems, University of Michigan, Ann Arbor, MI, USA

⁵Santa Fe Institute, Santa Fe, NM, USA

*Corresponding authors: E-mails: qianying@lanl.gov; eoromero@lanl.gov.

Associate editor: Adi Stern

Abstract

Reassortment is an evolutionary process common in viruses with segmented genomes. These viruses can swap whole genomic segments during cellular co-infection, giving rise to novel progeny formed from the mixture of parental segments. Since large-scale genome rearrangements have the potential to generate new phenotypes, reassortment is important to both evolutionary biology and public health research. However, statistical inference of the pattern of reassortment events from phylogenetic data is exceptionally difficult, potentially involving inference of general graphs in which individual segment trees are embedded. In this paper, we argue that, in general, the number and pattern of reassortment events are not identifiable from segment trees alone, even with theoretically ideal data. We call this fact the fundamental problem of reassortment, which we illustrate using the concept of the “first-infection tree,” a potentially counterfactual genealogy that would have been observed in the segment trees had no reassortment occurred. Further, we illustrate four additional problems that can arise logically in the inference of reassortment events and show, using simulated data, that these problems are not rare and can potentially distort our observation of reassortment even in small data sets. Finally, we discuss how existing methods can be augmented or adapted to account for not only the fundamental problem of reassortment, but also the four additional situations that can complicate the inference of reassortment.

Key words: phylodynamics, genomic reassortment, molecular epidemiology, population genetics.

Introduction

Reassortment is the evolutionary process by which viruses with segmented genomes exchange genetic material (Desselberger et al. 1978; Webster et al. 1992). When a cell becomes dually infected, that is, infected by two distinct viruses, the resulting viral offspring can contain a mixture of genomic segments from each parent (Simon-Loriere and Holmes 2011). The genetic diversity and fitness of segmented viruses such as influenza A virus (Macken et al. 2006; Rambaut et al. 2008), rotaviruses (Glass et al. 1994; McDonald et al. 2009), and bunyaviruses (Beaty et al. 1985; Briese et al. 2013) are influenced by reassortment. For example, in bunyaviruses, the M segment encodes the viral glycoprotein that mediates both cell fusion and immune responses (Elliott 2014). By swapping M segments, therefore, the reassortant can lead to much more severe illness in humans compared to its parents. As a concrete example, the hemorrhagic Ngari virus emerged from

reassortment between the only mildly pathogenic Bunyamwera and Batai viruses (Briese et al. 2006; Yanase et al. 2006). Other impacts of reassortment include escaping vaccine-elicited immunity (Batten et al. 2008), increasing fitness in humans (Matthijnsens et al. 2010), and getting access to new hosts (Martella et al. 2010). Understanding the molecular mechanism, the evolutionary process, and the epidemiological implications of reassortment is essential for anticipating and combating emerging infectious diseases (Morens and Fauci 2013), improving public health strategies (Morse 2001), and guiding the development of novel diagnostics, therapeutics, and vaccines (Rambaut et al. 2008). Further, reassortment can be thought of as a special case of recombination, one in which the breakpoints are fixed and correspond to the endpoints of each segment (Simon-Loriere and Holmes 2011).

In population genetics, a phylogenetic tree is commonly used to represent evolutionary history over time, in which

Received: September 20, 2023. **Revised:** February 23, 2024. **Accepted:** April 09, 2024

© The Author(s) 2024. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

Open Access

the root corresponds to the most recent common ancestor (MRCA) of all the sampled genomes, and internal nodes indicate a splitting of one lineage into two, denoting the ancestral relationships between sampled tips (Nixon 2001). For viral pathogens, this pattern of branching reflects both the epidemiology of who infected whom, and also exactly who within the population was sampled (Didelot et al. 2017). The logic underlying reconstructions of reassortment is that all else being equal, topological discordance between segment trees is caused by historical reassortment events (Fig. 1). For example, the observation that two viruses co-cluster in a tree inferred from one genomic segment, but are quite distant in a tree inferred from another segment, is likely the result of a recent reassortment event. This logic is the foundation of existing reassortment inference methods, which follow one of two distinct approaches. One approach begins with segment trees and then attempts to synthesize their topological incongruence into an ancestral reassortment graph (ARG), wherein specific reassortment events correspond to reticulations, *i.e.* partial merging of two ancestor lineages (Linder et al. 2004; Nagarajan and Kingsford 2010; Svinti et al. 2013; Barrat-Charlaix et al. 2022). This is analogous to the practice of reconciling gene trees to infer recombination (Maynard Smith and Smith 1998; Feil et al. 2001). These methods are intuitively appealing but they have trouble pinpointing the timing of reassortment events in dated segment trees (Collienne et al. 2024). A second approach seeks to infer an ARG—or a very closely related data structure called a phylogenetic network—jointly with the segment trees using the likelihood under a

coalescent-with-reassortment model (Müller et al. 2020; Stolz et al. 2022). These approaches are generally based on an extension of the Kingman coalescent (Kingman 1982; Wakeley 2005), which obtains a high level of computational efficiency by only considering the events in the population that are specifically ancestral to the sampled genomes. In the coalescent-based formulation, the likelihood of an ARG is computed backward in time (from tips to the root), where a coalescent event collapses two existing lineages into one and a reassortment event “creates” a new lineage from an existing one. This approach is analogous to the multi-species coalescent models used for joint gene tree–species tree inference (Degnan and Rosenberg 2009; Heled and Drummond 2009). However, both approaches share the core assumption that the best path to the inference of reassortment events is via an ARG-like object, itself derived from incongruities between the segment trees.

In this paper, we argue that reassortment inference faces a fundamental problem: the number and pattern of reassortment events is not identifiable without reference to some aspects of the underlying, forward-time, epidemiological process, namely, a latent data structure that we call the “first-infection tree.” We further identify four situations that we argue need to be explicitly handled by methods that infer reassortment. For this, we use a simple linear birth–death Markov chain simulation (Kendall 1948; Stadler 2010) to show that these situations are common even in small samples with perfect data for very simple epidemiological models. By simulating sequence data from these models, we also assess the error properties of

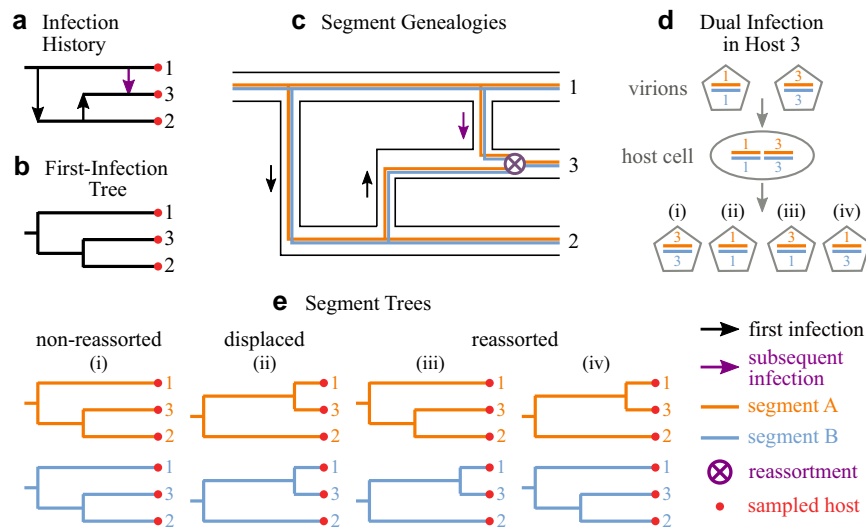


Fig. 1. Dual infection can lead to reassorted segment trees. Time flows from left to right in panels a), b), c), and e). a) The *infection history* encompasses all initial and subsequent infection events. In this small example, Host 1 infects Host 2, who then infects Host 3, and later Host 1 infects Host 3. Host 3 is then dually infected. b) The *first-infection tree* includes first infections of each host, but not subsequent infections or the identities of who infected whom. c) The genealogy of each of the two viral genomic segments can be traced within the infection history. d) Two virions in Host 3, the resident virion in Host 3 originated from the first infection from Host 2, and the invasive virion from the subsequent infection from Host 1, can reassort upon dual infection of a cell within Host 3. e) Depending on which virion is sampled from Host 3, the *segment trees* may show discordance with one another and/or with the first-infection tree. In outcome (i) the host maintains the viral genotype from its first infection, in (ii) the genotype from the second infection displaces the first, and in (iii) and (iv) a new reassorted viral genotype is generated.

commonly used methods, showing that estimation of the number of reassortment events necessary to explain incongruities in the reconstructed segment trees is not consistently correct. Our approach to both identifying problematic issues and designing a simulation model was to simplify biological reality as much as possible to obtain a core set of problems that are not dependent on particular model formulations or biological details. We have also attempted to design our forward simulation model to be as consistent as possible with existing methods using identical model forms and parameterizations where possible. This work leads us to a few general conclusions: (i) the inference of viral genomic reassortment is not separable from the inference of epidemiological dynamics, (ii) the ARG structure needs to be expanded to work for inference of viral genomic reassortment, and (iii) coalescent models are a difficult starting point for inferring reassortment. Our objective is not to criticize existing methods or suggest that previous work is wrong; in fact, we argue that there are several paths to robust methods that can be obtained from modifications of existing open-source approaches. Rather, our intention is to clarify the conceptual challenges in inferring reassortment for future implementations.

Results

Our results are laid out as follows. The first two sections illustrate logical problems that arise in the assessment of viral genomic reassortment. In those sections, we assume that all genealogies are binary and that reassortment is an instantaneous process that replaces one or more genomic segments in a given host with a copy from another extant host. In the final section, we use a forward linear birth–death simulation model to compute the frequency at which those problems arise in a theoretical study, and to simulate sequence and genealogical data to test methods for inferring viral reassortment. Details of the simulator and technical methods can be found in the Methods section and [Supplementary Material](#) online.

For clarity, we define a few standard and non-standard terms that we will use throughout the paper. We use the term *population history* to refer to the full history of all events that update the host dynamics or the genealogical structures in a population. In the context of a model, this is a record of who infected whom, when, and the times of all other sampling, removal, and reassortment events. [Figure 1](#) illustrates several important subsets of the population history. The *infection history* shows which individual infected (or re-infected) which other individual, and at what time. The *segment trees* are the observed genealogies of the individual viral segments in a host population. The *first-infection tree* is the subset of the infection history that only includes the first infection event (*i.e.* the creation of a newly infected individual) for each infected and sampled individual, without the specific information of who infected whom at the internal nodes. The first-infection tree can be thought of as the viral genealogy if no reassortment occurred, in which case the first-infection tree and

segment trees are all identical in topology and branch lengths. An infection history defines exactly one first-infection tree, but one first-infection tree is consistent with many possible infection histories. We assume that the true genealogy is directly observable, though in reality this involves inferring time-scaled phylogenetic segment trees. Also, we make no attempt to deal with the intricacies of tree reconstruction since the difficulties with incongruence-based reassortment inference can be made evident without doing so.

The Fundamental Problem

In this section, we illustrate the fundamental problem of reassortment using the hypothetical case where we have observed two segment trees and want to infer the true number of reassortment events required to explain their observed differences ([Fig. 2a](#)). Note that we make a distinction here between *visible* reassortment events that had some influence on the observed segment trees, and *invisible* reassortment events that are either not ancestral to the sampled viruses, or are ancestral but were replaced by a reassortment event that happened later in the same lineage ([Fig. 2b](#)).

A common way of counting the number of visible reassortment events is to find the minimum number of “remove-and-rejoin” operations ([Robinson 1971](#); [Svinti et al. 2013](#); [Collienne et al. 2024](#)) that are needed to reconcile the segment trees. The idea is to resolve any structural differences between segment trees by removing a tip or clade in one tree, along with its ancestral branch, and then rejoining that branch somewhere else in the same tree so as to reduce the differences between both trees. [Supplementary Fig. S3, Supplementary Material](#) online illustrates this idea. Generally, for such an approach the governing principle is parsimony: to propose the minimum number of reassortments that are needed to transform one segment tree into the other. Note that when we talk about incongruities in segment trees, we are referring to not only the tree topologies but also the branch lengths, which makes sense under the assumption that the segment trees can be observed without error. In practice, counting the number of reassortments might only rely on the inferred phylogenetic topology, under the assumption that differences in branch length are stochastic and/or due to evolutionary rate heterogeneity across tree branches. In that case, the remove-and-rejoin method is identical to the subtree-prune-regraft (SPR) method ([Robinson 1971](#)). However, since our goal is to identify the fundamental problem to inferring viral genomic reassortment, we assume that reassortment events which change the branch lengths while preserving the topology in the segment trees are visible.

The problem that arises in applying remove-and-rejoin moves on the segment trees is that, even when applied correctly, it does not necessarily find all the reassortment events that are visible in the data when one also considers the underlying infection history. [Figure 2a](#) illustrates the most basic case where the exact same set of segment trees

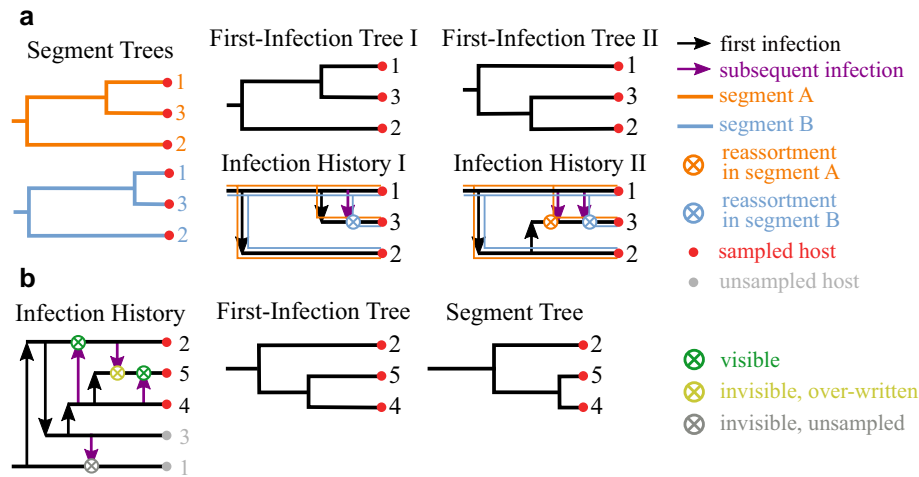


Fig. 2. How the infection history, first-infection tree, and sampling process impact genomic reassortment inference. a) Suppose that two segment trees are sampled. In a special case, First-Infection Tree I arose from Infection History I and is identical to one of the segment trees, so a single reassortment event is visible in the data. However, it could also be the case that First-Infection Tree II resulted from Infection History II, where two reassortment events are visible. Put another way, only one reassortment in segment B is needed to explain the segment trees' difference from First-Infection Tree I, while one reassortment in each segment is needed to explain why they both differ from First-Infection Tree II. b) Of all the reassortment events in the full population (all four ⊗ shown), only some leave a visible imprint in the data. Reassortment events that occur in unsampled lineages are invisible, as are those that occur in sampled lineages but are over-written by another reassortment event later in the same lineage.

can have a different number of visible reassortment events, depending on the infection history. While it is possible for a single remove-and-rejoin move (in this example) to convert one segment tree into the other, this resolution is not biologically possible under some infection histories. From this simple illustration, we obtain two important results: (i) reassortment events are not necessarily identifiable from segment trees without explicit reference to the unobserved population history, and therefore, (ii) inference of viral reassortment is fundamentally linked to inference of epidemiological dynamics. That is, to know even the number of visible reassortment events in the smallest possible tree, we need some way of accounting for the probabilities of all possible infection histories (e.g. an epidemiological model).

We argue that finding the number of visible reassortment events in the observed data does not actually require knowledge of the full infection history; rather, the first-infection tree alone is sufficient. The first-infection tree can be thought of as the tree structure that would have been observed in the segment trees if there had been no reassortment. It therefore represents a starting point for calculating the number of reassortment events required to obtain each of the segment trees (supplementary FigS3, Supplementary Material online). From this perspective, the minimum number of remove-and-rejoin moves is simply the minimum number of moves required to get from the first-infection tree to each of the segment trees. Starting from a different first-infection tree can give a different minimum number of reassortment events required to explain the data. Without specifying a first-infection tree as a starting point, remove-and-rejoin methods under parsimony implicitly assume that one of the segment trees has the same structure as the first-infection tree. For example, in Fig. 2a, unconditional remove-and-rejoin

parsimony will find that only one move is necessary, but this is correct only under the assumption that the first-infection tree (I) has the same structure as segment tree A. In general, failing to specify a first-infection tree will bias results to fewer reassortments because it makes the assumption that one of the segments has not reassorted. Therefore, unconditional parsimony is inconsistent as a criterion for reconstructing reassortment. Without either knowing (or assuming) the first-infection tree or assuming some model that gives an uncertainty expressed as a probability mass over the space of first-infection trees, parsimony cannot be a valid criterion for inferring reassortment. The broader issue also extends to methods that use softer constraints such as a penalty based on the norm of the reassortment rates or a prior in the context of a Bayesian analysis without explicitly conditioning on or sampling first-infection trees.

Four Situations That Make Inference of Reassortment Dynamics Difficult

In this section, we illustrate four situations that cause further problems for estimating quantities such as the reassortment rate and the time of reassortment events. We refer to these specific situations as “invisible,” “inaccurate,” “reversed,” and “obfuscated,” identifying where the potential for error arises in each case and the impact that it might have on reassortment inference. Any method that claims to be able to infer the history of reassortment in a population should be able to explicitly avoid or control for these errors.

Figure 3 illustrates the situation of *invisible* reassortment. This situation is the most straight-forward and easiest to correct. Invisible reassortment occurs when either a reassortment occurs in the part of the population that is unsampled, or when it is overwritten by another

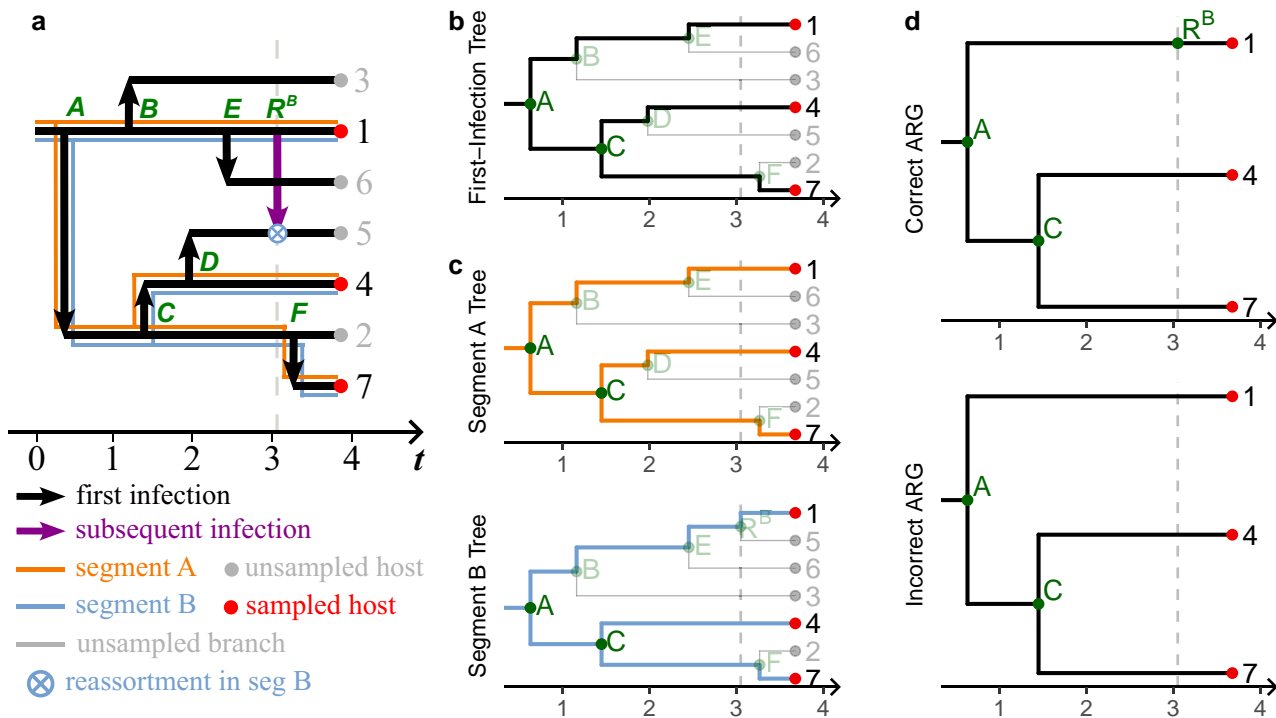


Fig. 3. Invisible reassortment. This figure is an example of an *invisible* type of reassortment event that arises when a reassortment event either occurs in an unsampled part of the population or is replaced by a later reassortment event in the sampled part of the population. Trees are measured in arbitrary time units increasing from left to right. The full infection history a), first-infection tree b), segment trees c), and ARGs d) are shown. The vertical grey dashed line indicates the time of that reassortment event.

reassortment event that happens later (see Fig. 2b). In the context of a structural analysis of reassortment, this is not viewed as problematic inasmuch as such analyses attempt to identify the minimum number of reassortments which explain the data and would not claim to detect invisible reassortments. However, for methods that attempt to infer the reassortment rate, failure to correct this situation will bias estimates of the reassortment rate. In Fig. 3, the invisible reassortment shown is not ancestral to the sample and, therefore, will only directly influence inferences concerning the broader effect of reassortment in the population (e.g. the distribution of fitness effects caused by any reassortment). However, invisible reassortment also occurs along sampled lineages (i.e. ancestral to the sample). Along a given lineage and its descendants, we can only observe the most recent successful reassortment event(s). Supplementary Fig. S2, Supplementary Material online illustrates how we dismiss those invisible reassortments along the lineages before performing the remove-and-rejoin moves to resolve the structural differences between trees. If a method does not correctly integrate the probability that an inferred ancestral reassortment event is visible, it could lead to an underestimation of the reassortment rate or spurious support for heterogeneous reassortment rates. This is a particularly acute problem for methods based on the coalescent, as they condition on the observed data and are unable to resolve this problem as the sampling rate in such models is not defined.

Figure 4 illustrates the situation of *inaccurate* inference of reassortment times that occurs when both individuals (or one of their direct descendants) involved in a reassortment are sampled. At time t_B , there is a reassortment event between samples 1 and 4, and, since both individuals are eventually sampled, the node corresponding to that event shows up as the most recent common ancestor of individuals 1 and 4 in the segment B tree. These types of nodes pose a problem for methods based on the manipulation of an ARG in the context of a reverse-time coalescent process. In the coalescent-with-reassortment model, the coalescent and reassortment processes must be assumed independent to ensure that the coalescent prior makes sense. However, as is clear in this example and from the representation of reassortment events in our model, an internal node in one or more of the sampled segment trees can be caused by a reassortment event rather than a first-infection event. In Fig. 4, the node R^B corresponds to a reassortment event that produced a node in segment tree B, where both children of that node (individuals 1 and 4) were sampled leading to the node associated with the reassortment event being in the segment B tree. To explain this, a coalescent-based approach needs to insert two nodes into the ARG, one for the reassortment event to split and another separate node for that additional lineage to coalesce. Since both nodes cannot co-occur, a bias is introduced. In principle, the distance between these two nodes can be made arbitrarily small such that they effectively occur at the same time. However, this is unlikely in

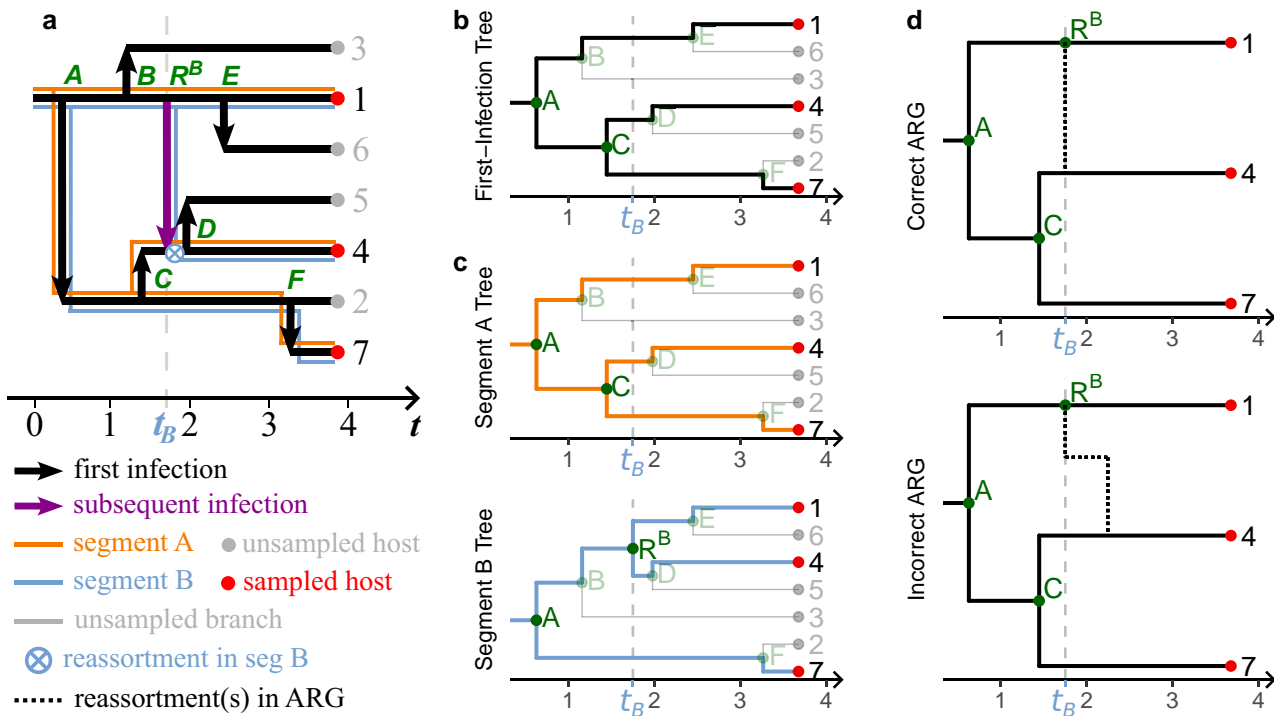


Fig. 4. Inaccurate reassortment. This figure is an example of an *inaccurate* reassortment event that arises when both children of a reassortment event are sampled creating a node in the segment tree that was caused by reassortment and not a new infection. The full infection history a), first-infection tree b), segment trees c), and ARGs d) are shown. The vertical gray dashed line indicates the time of that reassortment event and in d), the black dotted lines indicate an actual or inferred reassortment in the population.

practice unless the data are highly informative (*i.e.* very long, clean sequences). In the more common case where the data cannot entirely overcome the coalescent priors, the distance between the two nodes will be influenced by the coalescent model which, by definition, will have some non-zero mean, encouraging bias in the estimation of the reassortment and coalescent times.

Figure 5 illustrates the situation of *reversion* in the inference of reassortment. Reversion occurs when a segment is passed between sampled individuals, but only one offspring of the reassortment event is sampled, leaving the node corresponding to the reassortment event to be unobserved in the segment tree (*i.e.* individuals 5 and 6 are unsampled in Fig. 5c). The statistical situation of reversion has a similar effect to that of inaccuracy, in that either one or both of the timing of the reassortment and corresponding coalescent events in an ARG will not be correct. The cause is slightly different, however. When only one lineage descending from a reassortment event is sampled, the corresponding ARG requires a trifurcation (three lineages coming together in one node) at the sampled parent node of an unsampled reassortment event. Generally speaking, if some or all of the nodes of the ARG are assumed to follow a coalescent prior, then trifurcations will not be allowed, creating a similar situation as in the case of inaccuracy, where additional, spurious nodes must be inserted into the ARG (nodes A' and C' in Fig. 5d).

Figure 6 illustrates the situation of *obfuscation* in the inference of reassortment. The defining aspect of obfuscation

is that all segment trees have one or more reassortment events and at least one of these events is only partially sampled (*i.e.* only one of the descendants is sampled). We call this situation obfuscation because the first-infection tree no longer has the same structure as one of the segment trees. The effect of obfuscation is that the data appear to be explainable by fewer than the number of visible reassortments and the underlying infection history is thus completely obscured.

The Situations That Cause Problems for the Inference of Visible Reassortment Events are Common

Of the four situations that we discussed above, three are potential problems for inferring visible reassortments. In this section, we assess the frequency at which those situations arise theoretically and their impact on the accuracy of existing methods. We do not assess the frequency of invisible reassortment events as most reassortment events will be invisible, even in the context of an extremely large sample, and invisible reassortment events are not a problem in the assessment of visible reassortment events per se.

Using the population model detailed in the Methods section, we simulated a full epidemiological population history and extracted the sampled segment trees for a grid of different parameter combinations. We then counted the frequency of each of the different types of situations (details in the [Supplementary Material](#) online). Figure 7 shows the frequency of each situation for a range of simulation parameter values. In summary, the means

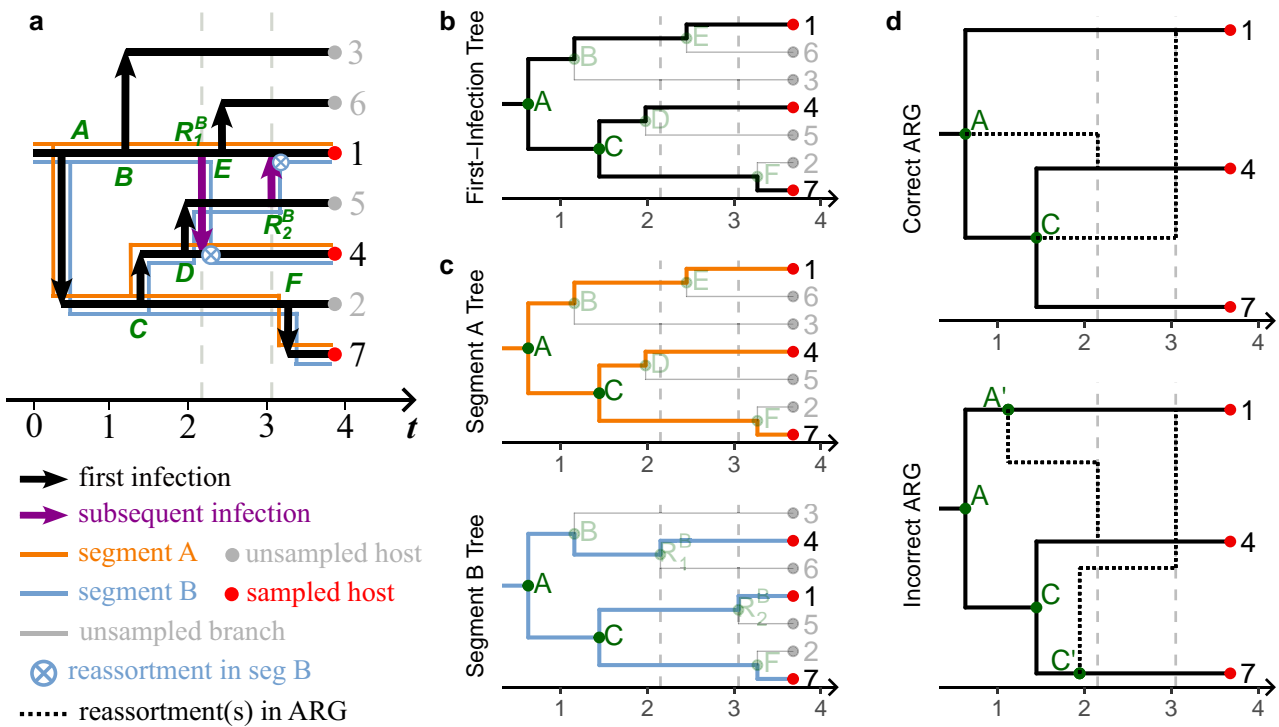


Fig. 5. Reversed reassortment. This figure is an example of a *reversed* reassortment event that arises when only one child of a reassortment event is sampled, creating the need for a trifurcating node in the ARG. The full infection history a), first-infection tree b), segment trees c), and ARGs d) are shown. The vertical gray dashed line indicates the time of that reassortment event and in d), the black dotted lines indicate an actual or inferred reassortment in the population.

of average frequencies for each situation are 26.1% (17%–45%), 6.7% (0.42%–9.6%), and 10.5% (0.49%–28.4%) for the situations *inaccuracy*, *reversion*, and *obfuscation*, respectively. The outcome variable is the proportion of nodes that have the type of referred situation averaged over 100 simulations. The *inaccuracy* situation increases as the first-infection rate becomes smaller than a fixed sampling rate. This makes sense as *inaccuracy* occurs when both descendants from a reassortment event are sampled, which is more likely when the first-infection rate is low, meaning that a larger proportion of nodes in the tree are caused by reassortment events, rather than first-infection events. Both *reversion* and *obfuscation* follow a similar pattern of being less frequent when the first-infection rate is high and the reassortment rate is low with respect to the sampling rate, and more frequent when the first-infection rate is low and the reassortment rate is high. Both of these types of situations occur when the nodes caused by reassortment events are not fully sampled (*i.e.* they do not appear as a node in the sampled segment trees), which will occur with higher frequency when a larger fraction of the nodes in a tree are caused by reassortment events rather than first-infection events.

Figure 8 shows the probability of inferring the number of visible reassortment events using two popular inference methods and an “omniscient” approach that counts the observed minimum number of remove-and-rejoin moves in 100 simulated data sets for each condition. In this context, the remove-and-rejoin approach is called omniscient

because we base it on the simulated population history rather than on inference, removing the issue of errors. In the case where reassortments only occur in one segment, the omniscient approach is always correct, by definition, but both the coalescent-based (Müller et al. 2020) and topology-based (Barrat-Charlaix et al. 2022) approaches tend to get the right number of reassortment events about 15% to 45% of the time. However, when both segments are allowed to reassort, the “omniscient” method also tends to get the number of reassortment events wrong, because one of the segment trees is no longer guaranteed to have the same structure as the first-infection tree. With the help of the simulated first-infection tree, the performances of both coalescent-based and topology-based methods are improved by about 10%. We further investigated the deviation in Coa1Re (Müller et al. 2020) estimations of reassortment number and the reassortment timing across all the posterior ARGs in [supplementary Figs. S4 and S5, Supplementary Material](#) online, respectively, showing that median errors in the assessment of reassortment times were generally about 12% to 18% of the total tree height (ranging from about 5% to 60%). Adding the simulated first-infection tree reduces the median errors to approximately 10%, ranging from 4% to 65%.

Discussion

We attempted to identify what we believe to be the core issues in the inference of viral genomic reassortment and

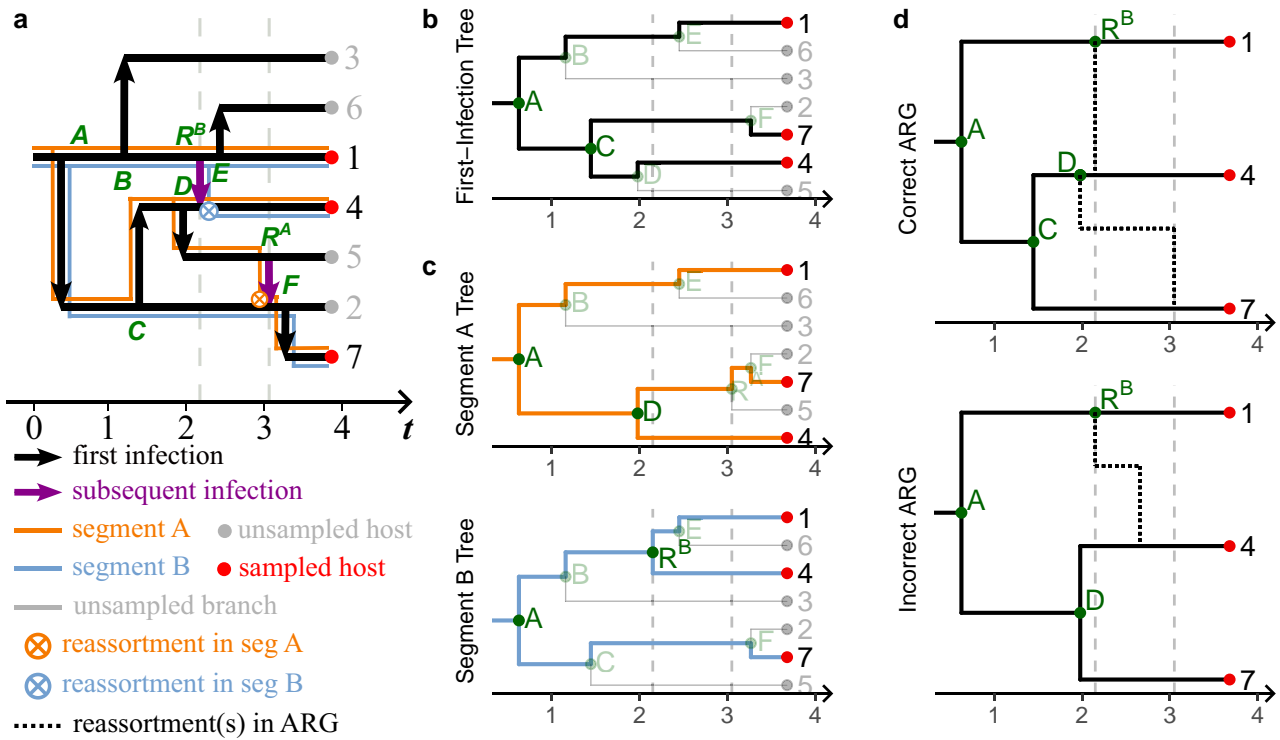


Fig. 6. Obscured reassortments. This figure is an example of an *obscured* reassortment event that arises when both sampled segment trees contain reassortment events such that the first-infection tree no longer has the same structure as one of the segment trees. The full infection history a), first-infection tree b), segment trees c), and ARGs d) are shown. The vertical gray dashed line indicates the time of that reassortment event and in d), the black dotted lines indicate an actual or inferred reassortment in the population.

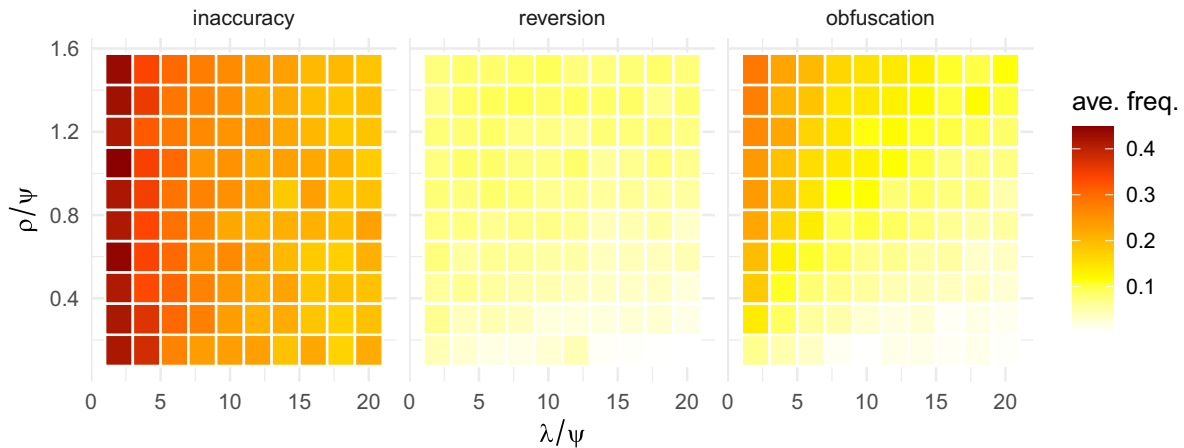


Fig. 7. The frequency of the *inaccuracy*, *reversion*, and the *obfuscation* situations. Starting from population size $l = 100$, the first infection rate is fixed at $\lambda = 5$ and the effective removal rate is $\mu + \psi = 5$, thus the expected population stays at $l = 100$. The total reassortment rate ρ is the sum of the rates for each segment and the joint reassortment rate: $\rho = \rho_A + \rho_B + \rho_{AB}$, where $\rho_A : \rho_B = 3 : 2$ and $\rho_{AB} = 0$. Samples are taken sequentially at a constant rate ψ until 50 individuals in the population are sampled.

argued that these issues are common and cause problems for existing approaches. Given the importance of viruses with segmented genomes to global public health, we see resolution of these issues as essential to achieve a more reliable understanding of the pandemic potential of those viruses.

Our analysis has important limitations. Elements of our results, such as the sufficiency of the first-infection

tree for resolving the inconsistency of parsimony, have not been proved or framed in a more rigorous context. While we believe these conclusions are likely to be proven true in the future, they should be understood to be contingent until more formal theories of the relationship between first-infection trees and the accuracy and precision of corresponding inferences can be established. Likewise, it is not so clear how the inconsistency of parsimony

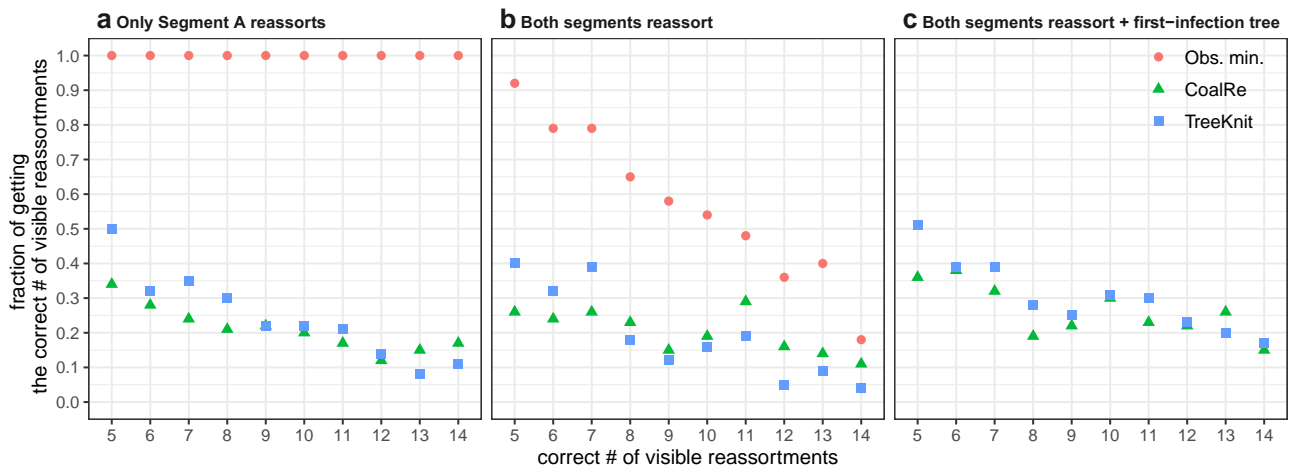


Fig. 8. Fraction of correctly estimating the number of visible reassortments. With a constant expected population size $l = 100$ and rates $\lambda = 5$, $\mu = 4.75$, $\psi = 0.25$, and $\rho = 0.2$, we simulated 100 sets of trees with 50 samples for each of the number of visible reassortments, ranging from 5 to 14. For each set, we computed the observed minimum number of reassortments and estimated the number of reassortments by the inferred MCC ARG from *CoalRe* and the ARG from *TreeKnit*, such that we computed the fraction of 100 simulated trees that estimate the number of visible reassortments correctly. Two reassortment schemes are used: in a), only segment A reassorts, $\rho_A = 0.2$ and $\rho_B = \rho_{AB} = 0$; in b) and c), both segments reassort but not simultaneously, $\rho_A = 0.12$, $\rho_B = 0.08$, and $\rho_{AB} = 0$. Two inference schemes are used: in a) and b), only segment trees are used for inference; in c), segment trees and the simulated first-infection tree are used.

for the inference of viral genomic reassortment affects either penalized maximum likelihood methods or Bayesian methods that do not explicitly infer or condition on a first-infection tree.

An important implication of this work is that the dynamics of infection and sampling in a population cannot be removed from the assessment of reassortment. That is, any claim about the history of reassortment implies a claim about the dynamics of infection that may be unreasonable or inconsistent with known facts. One of the ways that this makes reassortment a technically hard problem is that the coalescent framework that is popular in Bayesian phylogenetics (Bouckaert et al. 2019) and phylodynamics (Volz et al. 2009) is not an ideal starting point for inferring reassortment. First, the sampling rate in the coalescent is not defined as a population parameter. The coalescent conditions on the sample and works backward in time, such that the sampling rate (or probability in a batch sample) is not well defined. Even if the sampling rate can be inferred post-hoc by computing the implicit rate after-the-fact given the inferred population size dynamics (Parag et al. 2020), it would still be insufficient to account for invisible reassortment events because the coalescent is, by definition, silent on the population history where invisible reassortment events would most likely reside. Second, the coalescent has no obvious way of incorporating reassortment events in the context of the coalescent process itself. The coalescent has been proved to be an efficient approximation for inferring from a single genealogy, in terms of forward-time formulations such as the Wright-Fisher, Moran, and linear birth–death models (Kingman 1982; Wakeley 2005; King et al. 2022). Reassortment is a coupled death–birth event that generates a new branch in a segment tree without changing the population size, and there

is no theoretical foundation to link this forward-time model to the coalescent. Part of the reason that the coalescent is a popular formalism for Bayesian time-scaled phylogenetics is that it provides a theoretical basis for making otherwise weakly-identifiable evolutionary rate models work. However, in the context of reassortment, some or even most nodes might have nothing to do with the population process assumed by the coalescent, and, therefore, are either wrongly included in the coalescent prior or are left with an implicitly ill-defined prior. Further, the coalescent with reassortment will, by necessity, not allow trifurcations or a single node to be created by both a reassortment and coalescent event, which prevents such a method from addressing two of the four situations that we raised. Third, reassortment in a coalescent framework only makes sense when paired with a corresponding coalescent event (in reverse time, a reassortment node creates a new lineage in an ARG that will eventually need to coalesce), so both the reassortment rate and the coalescent process itself become uninterpretable. This is because the coalescent rate is, in most applications, computed as a function of the number of extant lineages and the effective population size. However, both first-infection and reassortment can result in a branching node in a tree, in which the former indeed increases population size, but the latter, formed by a coupled death–birth of an individual, does not change the population size. This convolution of the coalescent event and the reassortment event confounds the interpretation of their respective rates.

We argued that a reliable method for inferring reassortment both in the number and timing of events needs to address the fundamental problem with under-counting and the specific statistical situations that we raised. To our knowledge, no existing software or proposed method

meets these criteria, but there is very well-developed software that solves many of the most difficult aspects of the problem already. We see three practical ways that already-published methods could be adapted to address some or all of these issues.

We demonstrated that finding the most parsimonious number of remove-and-rejoin moves between sampled segment trees is not a consistent way to count the correct number of visible reassortment events, unless the first-infection tree is known. Existing topology-based methods (like those implemented in `TreeKnit` by [Barrat-Charlaix et al. 2022](#)) could be adapted to search the space of remove-and-rejoin moves for a set of segment trees given a first-infection tree. It is unlikely that the first-infection tree would be known, so a more reasonable way to deal with conditioning on the first-infection tree is to treat it as a nuisance variable that would be integrated out in the context of some prior distribution over the space of first-infection trees. Simulation methods could be used to generate first-infection trees from some model parameters, and existing remove-and-rejoin algorithms could be applied to all pairs of first-infection tree and segment tree. Effectively, this would both correct for the under-counting issue and would allow integration of prior knowledge in terms of the simulation model assumptions and the prior on the model parameters. Realistically, however, this approach might require a more theoretically robust understanding of how epidemiological dynamics constrain the inference of reassortment to make the computations feasible.

Methods that jointly search both the tree space and model parameter space provide not only a more complete reporting of uncertainty, but also direct estimates of the timing of reassortment events. The BEAST2 ([Bouckaert et al. 2019](#)) package `CoalRe` already implements most of the necessary functions to modify the process of sampling the joint posterior of first-infection trees and segment trees, rather than an ARG-like object in the context of a prior based on a forward-time model of epidemiological dynamics and sampling. Further, this approach could be adapted to model nonlinear epidemiological dynamics. We discuss the mathematical details of how this might work in the [Supplementary Material](#) online.

This paper's point could be summarized as, "to understand reassortment, one must place it in the context of a full population process." The problem is that there is no currently known method for computing the likelihood of either a sequence alignment or a set of segment trees given the population models that we consider in this paper. The issue lies in proposing the candidates of latent sampled first-infection trees and computing their likelihoods, which may be computationally challenging. Therefore, one potential way to deal with this is to treat the mapping between a population model and observable data through the lens of likelihood-free methods, such as Approximate Bayesian Computation (ABC) ([Tavaré et al. 1997](#)). These methods work by defining a metric space based on some

set of statistics computed on the observed data and data simulated from the population process of interest. Recently, [Voznica et al. \(2022\)](#) proposed a deep-learning approach along these lines where, rather than picking a set of statistics by hand, they allow a neural network to approximate a mapping between model parameters and simulated tree structures for population model selection, parameter estimation, and uncertainty quantification. A similar approach might prove fruitful for fitting complex models with genomic rearrangements such as reassortment.

A final consideration for the use of population models, as both explicit theories and data analytic tools, is that a well-defined population model can serve as a common domain for scientific hypotheses and the statistical tests needed to probe those hypotheses. Consider that genomic reassortment can occasionally produce highly fit viral strains that then rapidly spread in a population. It is easy to imagine how one might integrate this question into a population model, *e.g.* having a neutral type and an advantageous type (in terms of increased transmission rate) of reassortment. If we were able to fit that model to data, we could characterize a direct test of our scientific question as a clearly defined statistical test of the model (*e.g.* whether or not one transmission rate is larger than the other).

Methods

This section provides detailed interpretations and justifications on defining reassortment as a coupled death–birth event, incorporating reassortment in a simple population model with given parameterization, and simulating trees and sequences.

A Conceptual Operationalization of Reassortment in the Context of a Population Process

Reassortment is a complex process that involves cellular co-infection, competition among different viral strains, host immune responses, and the biophysics of genome packaging and virion formation. Each host-virus system will be unique in ways that may affect inferences regarding, for example, the timing of reassortment events. However, taking out the complex biology and within-host dynamics, the most basic operationalization of reassortment in a model is that of a coupled death–birth event. We illustrate how reassortment changes the evolutionary histories of viral segments in [supplementary Fig. S1, Supplementary Material](#) online for a hypothetical virus with two segments: A and B. In that case, host 0 co-infects host 2, swapping the B segment. Host 2's original B segment is lost in the population and their original branch in segment B is dissolved (*i.e.* the original branch "dies"). Correspondingly, host 2 now has a segment B that is directly descended from host 0 (*i.e.* a "birth" occurs in segment B at host 0). This type of coupled death–birth event is how reassortment produces conflicts in both topology and branch lengths in segment trees, due to differences in their evolutionary

histories (Vijaykrishna et al. 2015). Understanding the abstract effect of reassortment on genomic segment trees as this type of coupled death–birth event, in the context of a broader population process, is central to understanding the situations we raise in this paper.

This model defines a population process that consists of both the epidemiology of the host population and the ancestry of the viral population. Individuals can become infected and, once infected, can become infected again. We make the simplifying assumption that subsequent infection is resolved instantly, such that an individual only hosts one viral genotype at a time. We reserve the term *re-infection* for the situation in which the new viral genotype within the host differs from the previous one. To keep the model simple, we define reassortment as the combination of re-infection, cellular co-location, reassortment of one or more genomic segments, and dominance of the new reassortant at the host level. Each subset of genomic segments can, in principle, have a unique reassortment rate due to differences in within- or between-host survival, or the biophysical aspects of viral particle formation. Infected individuals leave the system by either dying or becoming sampled. For simplicity, we assume sampling is destructive to avoid the important but out-of-scope issues that arise in sampling direct descendants of already-sampled tips in phylogenetic trees. We assume all individuals are exchangeable, so there is no host- or virus-level heterogeneity, such as different behaviors or fitness.

Simulating the Frequency of Inaccuracy, Reversion, and Obfuscation

We implemented the conceptual model as a linear birth-death Markov process (Feller 1939; Kendall 1948) using the suggested formalism for simulating genealogical processes in King et al. (2022). The model contains two distinct levels of biological processes: epidemiology of the host population and ancestry within the viral population (supplementary table S1, Supplementary Material online). A conceptual description is provided here, and a more complete technical description is available in the Supplementary Material online and in King et al. (2022). The model assumes a viral infection with two segments (labeled A and B) and has six per-capita rates: the first-infection rate, λ ; the recovery rate, μ ; the sampling rate, ψ ; the reassortment rates for segment A and B, ρ_A and ρ_B , respectively; and the joint reassortment rate, ρ_{AB} , accounting for the simultaneous reassortment or the displacement as shown in Fig. 1e option (ii). Note that the effective removal rate of an individual is $\mu + \psi$, as we assume infected individuals who got sampled are immediately hospitalized or quarantined, and thus, removed from the infected population. The spread of infection in the host population is described by stochastic exponential growth at rate $\lambda - (\mu + \psi)$. The model records the full infection history (identity of source and recipient, and time of infection) in a format that enables much of the downstream computations that we use to count types of

reassortment inference situations. The first-infection tree is a subset of infection history that includes only the subtree of sampled tips and the nodes corresponding to the time they were first infected; source and recipient identity are also removed.

The descriptions above apply to the entire population of hosts and the viruses that infect them, *i.e.* the full population process. Not all individuals will be observed, however, and the sampled genealogy only represents the ancestral relationships between viral genomic samples. In the simulation, the full infection history, first-infection tree, and segment trees are pruned to include only the samples taken. Samples are taken sequentially from the infected population at the per-lineage sampling rate, ψ . For visual simplicity, figures illustrating stylized population histories or trees are always shown as if a homochronous, batch sample was taken at a given time.

All simulations start from one single infected individual ($l = 1$), while setting the first-infection rate $\lambda = 200$ and shutting down recovery, sampling, and reassortment until the population size reaches $l = 100$. We then turn on reassortment (*i.e.* $\rho_A > 0$ or/and $\rho_B > 0$, $\rho_{AB} = 0$) and run a burn-in period where the birth and removal rates are equal but there is no sampling (*i.e.* $\lambda = \mu = 5$ and $\psi = 0$) until the root time of the tree is later than the time of reaching $l = 100$, before proceeding with sampling (*i.e.* $\lambda = \mu + \psi$, where $\mu > 0$ and $\psi > 0$). We run the simulation until 50 samples are obtained. This approach guarantees that the expected population size stays constant at 100 during the studied period of the obtained trees. For each situation discussed in Results, we produced 100 stochastic replicates and computed the average frequency for each parameter combination.

Simulation and Analysis of Sequence Data

We first simulated the full population history such that the expected population size is ($l = 100$), with rates $\lambda = 5$, $\mu = 4.75$, $\psi = 0.25$, $\rho = \rho_A + \rho_B + \rho_{AB} = 0.2$, until 50 samples were obtained. These parameters are chosen to align with simulations in previous literature (Müller et al. 2020; Barrat-Charlaix et al. 2022), meanwhile making sure the frequencies of aforementioned situations are significant enough. For each of the number of visible reassortments, ranging from 5 to 14, we produced 100 replicates. Each segment genealogy represents precisely the history of that segment, but inferring a segment tree in practice requires genetic sequence data. Therefore, we used IQ-TREE 2 (Minh et al. 2020) to simulate alignments with a sequence length of 10^3 nucleotides using the Jukes–Cantor model for each tree, at an evolutionary rate 5×10^{-3} per site (Müller et al. 2020).

We analyzed our simulated data with two leading software packages for inferring reassortment. To use the CoalRe phylogenetic network model, we imported the alignments into BEAST2 and used the same mutation model. TreeKnit only requires the Newick format of tree data, so we directly input each set of trees.

Two different inference schemes were used: (i) using segment trees as input data only, and (ii) using the simulated first-infection tree and the segment trees as input data. In *CoalRe*, the simulated alignments for the first-infection tree were with a length of 10^4 nucleotides, approximately fixing its tree structure. In *TreeKnit*, we had pair-wise inputs, the first-infection tree and Segment A tree and the first-infection tree and Segment B tree, and summed up the number of reassortment events estimated separately from these inputs to have the total number of estimated reassortments.

Supplementary Material

Supplementary material is available at *Molecular Biology and Evolution* online.

Funding

This work was supported by grants from the U.S. National Institutes of Health (Grants #R01AI167048 to C.M.P., #R01AI087520 to T.L., #R01AI143852 to A.A.K.), and by the Joint U.S. National Science Foundation/National Institutes of Health Interface Program (Grant #1761603 to A.A.K.).

Conflict of Interest

None declared.

Data Availability

Code implementing the Markov Genealogical Processes with Reassortment and simulated data for figures are available at <https://github.com/MolEvolEpid/SimReassort>.

References

- Barrat-Charlaix P, Vaughan TG, Neher RA. TreeKnit: inferring ancestral reassortment graphs of influenza viruses. *PLoS Comput Biol*. 2022;**18**(8):e1010394. <https://doi.org/10.1371/journal.pcbi.1010394>.
- Batten CA, Maan S, Shaw AE, Maan NS, Mertens PPC. A European field strain of bluetongue virus derived from two parental vaccine strains by genome segment reassortment. *Virus Res*. 2008;**137**(1): 56–63. <https://doi.org/10.1016/j.virusres.2008.05.016>.
- Beaty BJ, Sundin DR, Chandler LJ, Bishop DHL. Evolution of bunyaviruses by genome reassortment in dually infected mosquitoes (*Aedes triseriatus*). *Science*. 1985;**230**(4725):548–550. <https://doi.org/10.1126/science.4048949>.
- Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, Heled J, Jones G, Kühnert D, De Maio N, et al. BEAST 2.5: an advanced software platform for Bayesian evolutionary analysis. *PLoS Comput Biol*. 2019;**15**(4):e1006650. <https://doi.org/10.1371/journal.pcbi.1006650>.
- Briese T, Bird B, Kapoor V, Nichol ST, Lipkin WI. Batai and Ngari viruses: M segment reassortment and association with severe febrile disease outbreaks in East Africa. *J Virol*. 2006;**80**(11): 5627–5630. <https://doi.org/10.1128/JVI.02448-05>.
- Briese T, Calisher CH, Higgs S. Viruses of the family Bunyaviridae: are all available isolates reassortants? *Virology*. 2013;**446**(1–2): 207–216. <https://doi.org/10.1016/j.virol.2013.07.030>.
- Collienne L, Whidden C, Gavryushkin A. Ranked subtree prune and regraft. *Bull Math Biol*. 2024;**86**(3):24. <https://doi.org/10.1007/s11538-023-01244-2>.
- Degnan JH, Rosenberg NA. Gene tree discordance, phylogenetic inference and the multispecies coalescent. *Trends Ecol Evol*. 2009;**24**(6):332–340. <https://doi.org/10.1016/j.tree.2009.01.009>.
- Desselberger U, Nakajima K, Alfino P, Pedersen FS, Haseltine WA, Hannoun C, Palese P. Biochemical evidence that “new” influenza virus strains in nature may arise by recombination (reassortment). *Proc Natl Acad Sci USA*. 1978;**75**(7):3341–3345. <https://doi.org/10.1073/pnas.75.7.3341>.
- Didelot X, Fraser C, Gardy J, Colijn C. Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Mol Biol Evol*. 2017;**34**(4):997–1007. <https://doi.org/10.1093/molbev/msw275>.
- Elliott RM. Orthobunyaviruses: recent genetic and structural insights. *Nat Rev Microbiol*. 2014;**12**(10):673–685. <https://doi.org/10.1038/nrmicro3332>.
- Feil EJ, Holmes EC, Bessen DE, Chan M-S, Day NPJ, Enright MC, Goldstein R, Hood DW, Kalia A, Moore CE, et al. Recombination within natural populations of pathogenic bacteria: short-term empirical estimates and long-term phylogenetic consequences. *Proc Natl Acad Sci USA*. 2001;**98**(1):182–187. <https://doi.org/10.1073/pnas.98.1.182>.
- Feller W. Die grundlagen der volterraschen theorie des kampfes ums dasein in wahrrscheinlichkeitstheoretischer behandlung. *Acta Biotheor*. 1939;**5**(1):11–40. <https://doi.org/10.1007/BF01602932>.
- Glass RI, Gentsch J, Smith JC. Rotavirus vaccines: success by reassortment? *Science*. 1994;**265**(5177):1389–1391. <https://doi.org/10.1126/science.8073280>.
- Heled J, Drummond AJ. Bayesian inference of species trees from multilocus data. *Mol Biol Evol*. 2009;**27**(3):570–580. <https://doi.org/10.1093/molbev/msp274>.
- Kendall DG. On the generalized “birth-and-death” process. *Ann Math Statist*. 1948;**19**(1):1–15. <https://doi.org/10.1214/aoms/1177730285>.
- King AA, Lin Q, Ionides EL. Markov genealogy processes. *Theor Popul Biol*. 2022;**143**:77–91. <https://doi.org/10.1016/j.tpb.2021.11.003>.
- Kingman JFC. The coalescent. *Stoch Process Their Appl*. 1982;**13**(3): 235–248. [https://doi.org/10.1016/0304-4149\(82\)90011-4](https://doi.org/10.1016/0304-4149(82)90011-4).
- Linder CR, Moret BME, Nakhleh L, Warnow T. Network (reticulate) evolution: biology, models and algorithms. In: Tutorial presented at the ninth Pacific symposium on biocomputing, 2004. <http://www.cs.rice.edu/~nakhleh/Papers/psb04.pdf>.
- Macken CA, Webby RJ, Bruno WJ. Genotype turnover by reassortment of replication complex genes from avian influenza A virus. *J Gen Virol*. 2006;**87**(10):2803–2815. <https://doi.org/10.1099/vir.0.81454-0>.
- Martella V, Bányai K, Matthijnsens J, Buonavoglia C, Ciarlet M. Zoonotic aspects of rotaviruses. *Vet Microbiol*. 2010;**140**(3–4): 246–255. <https://doi.org/10.1016/j.vetmic.2009.08.028>.
- Matthijnsens J, Rahman M, Ciarlet M, Zeller M, Heylen E, Nakagomi T, Uchida R, Hassan Z, Azim T, Nakagomi O, et al. Reassortment of human rotavirus gene segments into g11 rotavirus strains. *Emerging Infect Dis*. 2010;**16**(4):625–630. <https://doi.org/10.3201/eid1604.091591>.
- Maynard Smith J, Smith NH. Detecting recombination from gene trees. *Mol Biol Evol*. 1998;**15**(5):590–599. <https://doi.org/10.1093/oxfordjournals.molbev.a025960>.
- McDonald SM, Matthijnsens J, McAllen JK, Hine E, Overton L, Wang S, Lemey P, Zeller M, Van Ranst M, Spiro DJ, et al. Evolutionary dynamics of human rotaviruses: balancing reassortment with preferred genome constellations. *PLoS Pathog*. 2009;**5**(10): e1000634. <https://doi.org/10.1371/journal.ppat.1000634>.
- Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol*. 2020;**37**(5):1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
- Morens DM, Fauci AS. Emerging infectious diseases: threats to human health and global stability. *PLoS Pathog*. 2013;**9**(7): e1003467. <https://doi.org/10.1371/journal.ppat.1003467>.

- Morse SS. Factors in the emergence of infectious diseases. In: Price-Smith Andrew T., editor. *Plagues and politics: infectious disease and international policy*. New York: Palgrave; 2001. p. 8–26.
- Müller NF, Stolz U, Dudas G, Stadler T, Vaughan TG. Bayesian inference of reassortment networks reveals fitness benefits of reassortment in human influenza viruses. *Proc Natl Acad Sci USA*. 2020;**117**(29):17104–17111. <https://doi.org/10.1073/pnas.1918304117>.
- Nagarajan N, Kingsford C. GiRaF: robust, computational identification of influenza reassortments via graph mining. *Nucleic Acids Res*. 2010;**39**(6):e34. <https://doi.org/10.1093/nar/gkq1232>.
- Nixon KC. Phylogeny. In: Levin SA, editor. *Encyclopedia of biodiversity*. 2nd ed. Amsterdam: Elsevier; 2001. p. 16–23.
- Parag KV, du Plessis L, Pybus OG. Jointly inferring the dynamics of population size and sampling intensity from molecular sequences. *Mol Biol Evol*. 2020;**37**(8):2414–2429. <https://doi.org/10.1093/molbev/msaa016>.
- Rambaut A, Pybus OG, Nelson MI, Viboud C, Taubenberger JK, Holmes EC. The genomic and epidemiological dynamics of human influenza A virus. *Nature*. 2008;**453**(7195):615–619. <https://doi.org/10.1038/nature06945>.
- Robinson DF. Comparison of labeled trees with valency three. *J Comb Theory Ser B*. 1971;**11**(2):105–119. [https://doi.org/10.1016/0095-8956\(71\)90020-7](https://doi.org/10.1016/0095-8956(71)90020-7).
- Simon-Loriere E, Holmes EC. Why do RNA viruses recombine? *Nat Rev Microbiol*. 2011;**9**(8):617–626. <https://doi.org/10.1038/nrmicro2614>.
- Stadler T. Sampling-through-time in birth-death trees. *J Theor Biol*. 2010;**267**(3):396–404. <https://doi.org/10.1016/j.jtbi.2010.09.010>.
- Stolz U, Stadler T, Müller NF, Vaughan TG. Joint inference of migration and reassortment patterns for viruses with segmented genomes. *Mol Biol Evol*. 2022;**39**(1):msab342. <https://doi.org/10.1093/molbev/msab342>.
- Svinti V, Cotton JA, McInerney JO. New approaches for unravelling reassortment pathways. *BMC Evol Biol*. 2013;**13**(1):1–14. <https://doi.org/10.1186/1471-2148-13-1>.
- Tavaré S, Balding DJ, Griffiths RC, Donnelly P. Inferring coalescence times from DNA sequence data. *Genetics*. 1997;**145**(2):505–518. <https://doi.org/10.1093/genetics/145.2.505>.
- Vijaykrishna D, Mukerji R, Smith GJD. RNA virus reassortment: an evolutionary mechanism for host jumps and immune evasion. *PLoS Pathog*. 2015;**11**(7):e1004902. <https://doi.org/10.1371/journal.ppat.1004902>.
- Volz EM, Kosakovsky Pond SL, Ward MJ, Leigh Brown AJ, Frost SDW. Phylodynamics of infectious disease epidemics. *Genetics*. 2009;**183**(4):1421–1430. <https://doi.org/10.1534/genetics.109.106021>.
- Voznica J, Zhukova A, Boskova V, Saulnier E, Lemoine F, Moslonka-Lefebvre M, Gascuel O. Deep learning from phylogenies to uncover the epidemiological dynamics of outbreaks. *Nat Commun*. 2022;**13**(1):3896. <https://doi.org/10.1038/s41467-022-31511-0>.
- Wakeley J. *Coalescent theory, an introduction*. Greenwood Village (CO): Roberts and Company; 2005.
- Webster RG, Bean WJ, Gorman OT, Chambers TM, Kawaoka Y. Evolution and ecology of influenza A viruses. *Microbiol Rev*. 1992;**56**(1):152–179. <https://doi.org/10.1128/mr.56.1.152-179.1992>.
- Yanase T, Kato T, Yamakawa M, Takayoshi K, Nakamura K, Kokuba T, Tsuda T. Genetic characterization of Batai virus indicates a genomic reassortment between orthobunyaviruses in nature. *Arch Virol*. 2006;**151**(11):2253–2260. <https://doi.org/10.1007/s00705-006-0808-x>.