

# QUANTIFYING GENETIC INNOVATION: MATHEMATICAL FOUNDATIONS FOR THE TOPOLOGICAL STUDY OF RETICULATE EVOLUTION

MICHAEL LESNICK, RAÚL RABADÁN, AND DANIEL I. S. ROSENBLOOM

**ABSTRACT.** A topological approach to the study of genetic recombination, based on persistent homology, was introduced by Chan, Carlsson, and Rabadán in 2013. This associates a sequence of signatures called *barcodes* to genomic data sampled from an evolutionary history. In this paper, we develop theoretical foundations for this approach.

First, we clarify the motivation for the topological approach by presenting a novel formulation of the underlying inference problem. Specifically, we introduce and study the *novelty profile*, a simple, stable statistic of an evolutionary history which not only counts recombination events but also quantifies how recombination creates genetic diversity. We propose that the (hitherto implicit) goal of the topological approach to recombination is the estimation of novelty profiles.

We then develop a theory for the estimation of novelty profiles using barcodes. We focus on a low-recombination regime, where the evolutionary history can be described by a directed acyclic graph called a *galled tree*, which differs from a tree only by isolated topological defects. We show that in this regime, under a complete sampling assumption, the 1<sup>st</sup> barcode yields a lower bound on the novelty profile, and hence on the number of recombination events. For  $i > 1$ , the  $i^{\text{th}}$  barcode is empty. In addition, we use a stability principle to strengthen these results to ones which hold for any subsample of an arbitrary evolutionary history. To establish these results, we describe the topology of the Vietoris-Rips filtrations arising from evolutionary histories indexed by galled trees.

As a step towards a probabilistic theory, we also show that for a random history indexed by a fixed galled tree and satisfying biologically reasonable conditions, the intervals of the 1<sup>st</sup> barcode are independent random variables. Using simulations, we explore the sensitivity of these intervals to recombination.

## CONTENTS

1. Introduction	2
Acknowledgements	6
2. Phylogenetic Graphs and Evolutionary Histories	6
3. Novelty Profiles	8
3.1. The Temporal Novelty Profile	8
3.2. The Topological Novelty Profile	9
4. Histories Indexed By Galled Trees	11
5. Topological Preliminaries	13
5.1. Persistent Homology	13
5.2. Stability of the Topological Novelty Profile	16
5.3. Discrete Morse Theory	16
6. Barcodes of Histories Indexed by Galled Trees	17

6.1. Barcodes of Histories indexed by Trees	17
6.2. Metric Decomposition of an Evolutionary History	18
6.3. Vietoris-Rips Filtrations of Almost Linear Metric Spaces	19
6.4. Inference about Recombination from Barcodes	24
7. Relaxing the Complete Sampling and Galled Tree Assumptions	26
7.1. Relaxing the Complete Sampling Assumption	26
7.2. Relaxing the Galled Tree Assumption	26
8. Random Histories Indexed by Galled Trees	27
8.1. Independence of Intervals in the First Barcode	27
8.2. The Barcode of a Random History on a Directed Loop: Numerical Results	31
9. Discussion	34
Appendix A. Probability that the Coalescent with Recombination Generates a Galled Tree	37
A.1. Overview of the Coalescent with Recombination	37
A.2. Probability of Generating a Galled Tree as the Solution to a Linear System	37
Appendix B. Subsamples Rarely Violate our Theoretical Bounds for Complete Samples	40
References	41

## 1. INTRODUCTION

Recombination is a process by which the genomes of two parental organisms combine to form a new genome. Like genetic mutation, recombination gives rise to genetic diversity in evolving populations. But unlike mutation, recombination can unite advantageous traits which have arisen in separate lineages, or rescue an advantageous trait from an otherwise disadvantageous genetic background. In these ways, recombination hastens the pace at which beneficial genetic novelty arises.

Evolving populations can be studied by observing genetic sequences obtained from a sample of organisms. Several methods exist to estimate the number of recombination events that have occurred in the ancestry of a sample and to identify the genomic locations where recombination may have occurred [32, 41, 46]. Yet these methods do not reveal how recombination generates genetic diversity: Recombination between two very distinct parents may create a genetically very novel offspring, contributing substantial diversity to the population, but recombination between genetically similar parents can only create genetically similar offspring, contributing little diversity.

In this work, we introduce a simple, stable statistic of an evolving population, the *novelty profile*, which quantifies how recombination contributes to genetic diversity. To define the novelty profile, we first need to select a formal model of an evolving population. We call the model we consider in this paper an *evolutionary history*. An evolutionary history  $E$  is a directed acyclic graph  $G$ , together with a set  $E_v$  at each vertex  $v$  of  $G$ , satisfying certain conditions. We call  $G$  a *phylogenetic graph*, and say that  $G$  *indexes*  $E$ . Each vertex of  $G$  represents an organism, each edge of  $G$  represents a parental relationship, and each  $E_v$  specifies the genome of the organism  $v$ . See Section 2 for the formal definition of an evolutionary history and an illustration.

The novelty profile of an evolutionary history is simply a list of  $k$  monotonically decreasing numbers, where  $k$  is the number of recombination events in the history. Roughly, each number measures the contribution to genetic diversity of one recombinant. We introduce two versions of this statistic, the *temporal* and *topological* novelty profiles. The definition of the temporal novelty profile is very elementary and intuitive, but depends on a specification of the time at which each organism is born. Moreover, the temporal novelty profile, while stable to perturbations (i.e., small changes) of the genomes, is unstable to perturbations of the birth times. In contrast, the topological novelty profile is defined in a way that does not depend on birth times. It is also stable to perturbations of the genomes. The topological novelty profile is a lower bound for the temporal novelty profile, in the sense that the  $i^{\text{th}}$  element of the topological novelty profile is less than or equal to the  $i^{\text{th}}$  element of the temporal novelty profile for all  $i$ .

We consider the estimation of novelty profiles using *Persistent homology*, a popular topological data analysis method. The topological approach to recombination has previously been studied in a series of papers [8, 9, 12, 23, 24]. In this previous work, persistent homology was used to associate a sequence  $\mathcal{B}_0(S), \mathcal{B}_1(S), \mathcal{B}_2(S), \dots$  of objects called *barcodes* to an arbitrary sample  $S$  of an evolving population. Each barcode is a collection of intervals  $[a, b]$  on the real line. In [12], it was shown that, under a standard *infinite sites* assumption ruling out multiple mutations at the same genetic site, if no recombination occurs in the population's history, then  $\mathcal{B}_i(S)$  is empty for all  $i \geq 1$ . Hence a non-empty barcode  $\mathcal{B}_i(S)$  for any  $i \geq 1$  certifies that recombination has occurred at some point in the history.

Within simulations of evolving populations, the number of intervals in the first barcode  $\mathcal{B}_1(S)$  has been observed to increase with the simulated recombination rate [9]. It has also been observed that the endpoints of the intervals depend on certain aspects of population structure and genetic diversity [8, 23]. However, the precise statistical nature of these relationships has not been made clear, nor is it clear in what generality these relationships hold. Our work is inspired by and aims to explain these observations.

We propose that the central inference problem implicit in the previous topological work on evolution is the estimation of novelty profiles. We study barcodes as estimators of the novelty profile, focusing attention on a low-recombination regime where the evolutionary histories are indexed by *galled trees*. Galled trees are directed acyclic graphs that are almost trees, in a sense: They may have cycles, but these cycles are topologically separated from one another; see Section 4 for the precise definition. Galled trees have received considerable attention in the phylogenetics literature as computationally convenient models of evolution with infrequent recombination [28, 33]. They have been of interest primarily because certain phylogenetic network reconstruction problems that are computationally hard in general admit polynomial-time solutions when restricted to galled trees. To clarify the biological relevance of galled tree models of evolution, in Appendix A we study the probability  $P$  that a galled tree correctly models an evolutionary history. We work with a coalescent model of evolution, a standard model in population genetics. We show that for this model,  $P$  can be computed by solving a linear system of equations, and we observe that for a fixed population size,  $P$  tends to 1 as the recombination rate tends to 0; see also Remark 4.6.

We observe that for evolutionary histories indexed by galled trees, the temporal and topological novelty profiles are equal (Proposition 4.5). Our main result relating barcodes to

recombination in the galled tree setting is the following (see Theorem 6.18 and the preceding definitions for the precise formulation):

**Theorem.** *Let  $\mathcal{E}$  be an evolutionary history indexed by a galled tree.*

- (i) *The set of lengths of intervals in the barcode  $\mathcal{B}_1(\mathcal{E})$  is a lower bound on the novelty profile. In particular, the number of intervals in  $\mathcal{B}_1(\mathcal{E})$  is a lower bound on the number of recombination events in  $\mathcal{E}$ .*
- (ii)  *$\mathcal{B}_i(\mathcal{E})$  is empty for  $i \geq 2$ .*

Part (i) of the theorem does not hold for barcodes  $\mathcal{B}_1(S)$  of arbitrary samples  $S \subset \mathcal{E}$  (Example 6.19). However, using a well-known stability property of persistent homology, we observe that the theorem extends to an approximate version which holds for an arbitrary sample  $S$ , even in the presence of noise (Corollary 7.1). The quality of the approximation depends on the similarity of the geometries of  $S$  and  $\mathcal{E}$ , as measured by the *Gromov-Hausdorff distance*. Along similar lines, the theorem further extends to an approximate version for histories indexed by arbitrary phylogenetic graphs (Corollary 7.3); here, the quality of approximation is controlled by the number of mutations which must be ignored to obtain a history indexed by a galled tree.

These results are deterministic; in cases where the history is sampled at random from a known distribution, one hopes to be able to obtain stronger probabilistic results. As a first step towards such results, we show in Section 8.1 that for a random history indexed by a fixed galled tree and satisfying a biologically reasonable independence condition, the intervals of the 1<sup>st</sup> barcode are independent random variables indexed by the recombinants of  $G$ .

We then study the distributions of these random variables via simulation, for one class of random models of genetic sequence evolution. Our simulation results indicate that even when we have sampled all individuals in the evolutionary history, the barcode often substantially underestimates the novelty profile. For example, in the most favorable circumstances, a recombinant of high novelty is detected in our simulations about a third of the time. Nevertheless, the barcodes allow us to deduce partial information about the novelty profile. Notably, we observe in our simulations that when a recombination event of novelty  $n$  is detected by the barcode, the average length of the corresponding interval is approximately  $c + d\sqrt{n}$  for constants  $c$  and  $d$ .

*The related work of Cámara et al.* Theoretical foundations for the application of persistent homology to recombination have also been studied in recent work of Cámara, Levine, Rabadán [8], though from a rather different angle than ours. That work considers connections to the problem of constructing *minimal ancestral recombination graphs* (ARGs) for single-breakpoint models of recombination. (An ARG roughly corresponds to what we call an evolutionary history in this paper; a minimal ARG is a history of a given set of genome sequences with as few recombination events as possible.) In contrast, we do not consider ARG reconstruction or constrain recombination to a single-breakpoint model. While our aim and technical approach differ from this previous work, we do share the common goal of understanding topological statistics for evolutionary biology.

*Mathematical Contributions.* One key feature of the barcode signatures of recombination studied here is that they depend only on the *metric structure* on an evolutionary history,

i.e., the genetic distances between organisms—in our formalism, the Hamming distance, or monotonic transformations thereof such as the Jukes-Cantor distance. In fact, these barcodes are given by a standard construction which associates barcodes to any finite metric space  $M$ . In this construction, one first builds a 1-parameter family of simplicial complexes  $\mathcal{V}(M)$  called the *Vietoris-Rips filtration (VRF)*.

The topological study of VRFs is a central theoretical problem in topological data analysis. While some fundamental results about VRFs are well known, including a stability theorem [7, 14, 16], relatively little is known about concrete computations of the topology of VRFs, outside of special cases; even for points distributed uniformly on a circle or ellipse, the problem is already non-trivial, and has been the subject of recent research [1, 3].

The result of Chen et al., that  $\mathcal{B}_i(S) = \emptyset$  for  $i \geq 1$  when  $S$  is a sample of a history with no recombination, amounts to a proof that the VRF of a *tree-like* metric space is topologically trivial, up to multiplicity of connected components; see [12] or Proposition 6.3. Analogously, the mathematical heart of our main results about recombination for galled trees is a topological description of  $\mathcal{V}(\mathcal{E})$ , for  $\mathcal{E}$  an evolutionary history indexed by a galled tree, regarded as a metric space: We use discrete Morse theory [26] to show that each simplicial complex in  $\mathcal{V}(\mathcal{E})$  is homotopy equivalent to a disjoint union of bouquets of circles, where each circle corresponds to a unique recombination event. Moreover, we completely describe the topological behavior of the inclusion maps in  $\mathcal{V}(\mathcal{E})$  and give bounds on the number of the intervals in  $\mathcal{B}_1(\mathcal{E})$ . For the precise statements, see Proposition 6.12 and Theorems 6.11 and 6.13.

Our topological study of  $\mathcal{V}(\mathcal{E})$  hinges on the study of the VRFs of *almost linear metric spaces*; we say a metric space is almost linear if (up to isometry) it is obtained from a finite subspace of  $\mathbb{R}$  by adding a single point. In brief, almost linear metric spaces enter into our analysis in the following way: We observe in Proposition 6.10 that, up to isometry, the metric space  $\mathcal{E}$  can be constructed by iteratively gluing together tree-like metric spaces and almost linear metric spaces using a coproduct construction. (The coproducts are taken in a category of based metric spaces, allowing the basepoint to change.) Moreover, letting  $P \vee Q$  denote a coproduct of two based metric spaces  $P$  and  $Q$ , we have that  $\mathcal{V}(P \vee Q)$  is, up to homotopy, a wedge sum of  $\mathcal{V}(P)$  and  $\mathcal{V}(Q)$  (Proposition 6.7). Since the VRFs of tree-like metric spaces are topologically trivial, it follows that to describe the topology of  $\mathcal{V}(\mathcal{E})$ , it suffices to describe the topology of the VRF of an almost linear metric space; Theorem 6.13 gives such a description.

*Outline.* For some of the material of this paper, we must assume that the reader is familiar with elementary algebraic topology. However, much of our material on novelty profiles does not require a background in topology, and we believe that this material may be of interest, independent of the topological material. Thus, we have arranged the paper so that the material on topology appears as late as possible.

Section 2 introduces our mathematical formalism for working with evolving populations in the presence of recombination. Section 3 introduces novelty profiles. Section 4 reviews galled trees and establishes that in the special case of galled trees, the temporal and topological novelty profiles are equal. Section 5 reviews aspects of persistent homology and discrete Morse theory needed in the remainder of our paper, and observes that the topological novelty profile is stable. Section 6.1 briefly reviews the results of Chan et al. on barcodes of evolutionary

histories indexed by trees. Section 6.2 studies the VRFs of coproducts of based metric spaces, and Section 6.3 presents our topological analysis of the VRFs of almost linear metric spaces. Using the results of Sections 6.1 to 6.3, Section 6.4 establishes our main deterministic result about the barcodes of evolutionary histories indexed by galled trees. Section 7 applies the stability of persistent homology to extend this result to subsamples of histories indexed by arbitrary phylogenetic graphs.

Section 8.1 establishes that for a suitably chosen random history indexed by a fixed galled tree, the intervals in the 1<sup>st</sup> persistence barcode are independent random variables. With this as motivation, Section 8.2 uses simulation to study the statistical properties of the barcode of a random history with a single recombination event. Section 9 discusses the applicability of our results to real-world genomic data, and explores directions for future work.

Two appendices tie our results explicitly to coalescent theory. Appendix A studies the probability that an evolutionary history generated by the coalescent model is a galled tree. Appendix B observes in simulation that, although our main result for histories indexed by galled trees does not hold exactly for arbitrary subsamples, the lower bound on the number of recombination events implied by that result is only rarely violated under subsampling.

**Acknowledgements.** We thank Ulrich Bauer for helpful discussions about how to prove Theorem 6.13, our main result about the Vietoris-Rips filtrations of almost linear metric spaces. Ulrich provided valuable input about the use of the triangle inequality in that argument, and also suggested the use of the discrete gradient vector field of [35]. We also thank Pablo Cámara and Kevin Emmett for valuable discussions, and we thank Greg Henselman for helpful feedback on our discussion of discrete Morse theory. Lesnick was partially supported by funding from the Institute for Mathematics and its Applications, NIH grants U54CA193313 and T32MH065214, and an award from the J. Insley Blair Pyne Fund. Rabadán and Rosenbloom were funded by NIH grants U54CA193313 and R01GM117591.

## 2. PHYLOGENETIC GRAPHS AND EVOLUTIONARY HISTORIES

We now introduce our mathematical formalism for the topological study of reticulate evolution. The formalism is similar to that used elsewhere in the literature on reticulate evolution, though some of our terminology is non-standard; for context, see for example [33] and the references therein.

**Definition 2.1** (Phylogenetic Graph). A *phylogenetic graph* is a finite directed acyclic graph  $G$  such that

1.  $G$  has a unique vertex, the *root*, with in-degree 0,
2. Each vertex of  $G$  has in-degree at most 2.

We call a vertex in  $G$  of in-degree 1 a *clone*, and a vertex of in-degree 2 a *recombinant*. If  $(v, w)$  is a directed edge in  $G$ , we say  $v$  is a *parent* of  $w$ . We define a *rooted tree* to be a phylogenetic graph with no recombinants.

Fig. 1 illustrates a simple phylogenetic graph.

For  $G$  a rooted directed acyclic graph with vertex set  $V$  and  $S \subset V$ , we say  $v \in S$  is the *minimum* of  $S$  if for all  $s \in S$ , any directed path from  $r$  to  $s$  in  $G$  contains  $v$ .  $S$  may not have a minimum element, but if the minimum element exists, it is clearly unique.

Let **Set** denote the collection of all finite sets.



**Definition 2.5** (Symmetric Difference Metric). Define a metric  $d$  on finite sets, the *symmetric difference metric* by taking

$$d(A, B) = |(A \cup B) \setminus (A \cap B)|$$

for any finite sets  $A, B$ . For any history  $\mathcal{E}$  indexed by a phylogenetic graph  $G$  with vertex set  $V$ , this restricts to a metric on the set  $\{\mathcal{E}_v \mid v \in V\}$ . We denote the resulting metric space as  $\text{met } \mathcal{E}$ , or when no confusion is likely, simply as  $\mathcal{E}$ .

**Remark 2.6.** It is common in the phylogenetics literature to model genomes as binary vectors, and to metrize a set of genomes using the Hamming distance. It is easy to see that the formalism we’ve introduced here is essentially equivalent. Under this equivalence, other common phylogenetic distances (e.g., Jukes-Cantor distance, Nei-Tamura distance) correspond to monotonic transformations of the symmetric difference metric  $d$ . In fact, all the results of this paper formulated in terms of  $d$  extend immediately to such monotonic transformations.

**Remark 2.7.** In real-world evolving populations, the infinite sites assumption, described above, may not always hold. In other words, the same mutation may occur in different organisms despite being absent in their common ancestors. Such mutations, termed *homoplasies*, may be observed in sampled data either if the per-site mutation rate is high (which is typical for species with short genomes, such as RNA viruses) or if the mutations confer high fitness. Homoplasies are typically rare for species with long genomes, as the probability of mutating twice at the same exact genetic site is small. If they do occur, homoplasies usually involve few sites, so that the metric space underlying the history differs only slightly from that of a history satisfying the infinite sites assumption.

### 3. NOVELTY PROFILES

**3.1. The Temporal Novelty Profile.** For  $G$  a phylogenetic graph with vertex set  $V$ , say  $t : V \rightarrow \mathbb{R}$  is a *time function* if  $t(v) < t(w)$  whenever  $v < w$  in the partial order on  $V$  induced by  $G$ . We interpret  $t(v)$  as the birth time of organism  $v$ .

**Definition 3.1** (Temporal Novelty Profile). Given a history  $\mathcal{E}$  indexed by  $G$ , a time function  $t : V \rightarrow \mathbb{R}$ , and a recombinant  $r$  of  $G$ , we define  $\mathcal{N}(r)$ , the temporal novelty of  $r$ , by

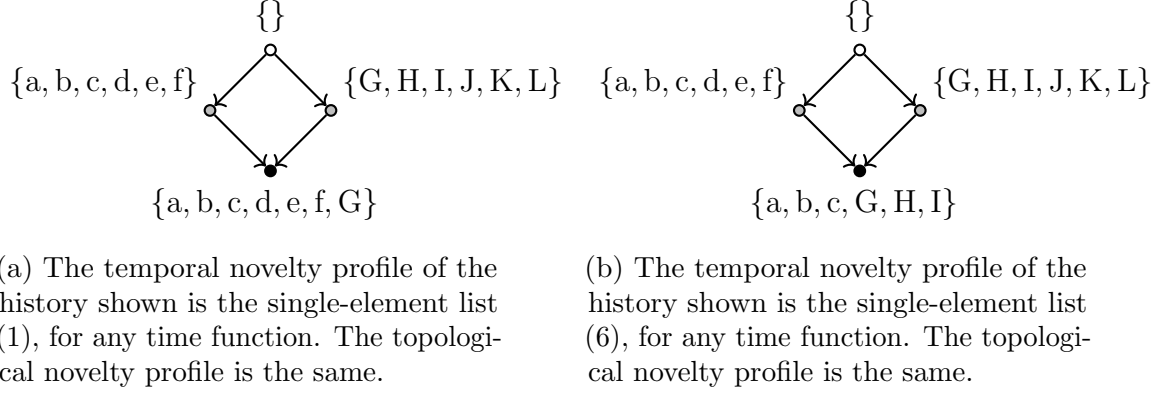
$$\mathcal{N}(r) := \min \{d(\mathcal{E}_v, \mathcal{E}_r) \mid t(v) < t(r)\}.$$

We define  $\mathcal{N}(\mathcal{E}, t)$ , the *temporal novelty profile* of  $\mathcal{E}$  (with respect to  $t$ ) to be the list of temporal novelties  $\mathcal{N}(r)$ , for all recombinants  $r$  of  $G$ , sorted in decreasing order.

**Example 3.2.** Fig. 3 illustrates novelty profiles for two histories indexed by the same simple phylogenetic graph. For any time function on the history shown in Fig. 3a, the unique recombinant has temporal novelty 1, so the temporal novelty profile of this history is the single-element list (1). Similarly, for any time function on the history shown in Fig. 3b, the temporal novelty profile is the single-element list (6).

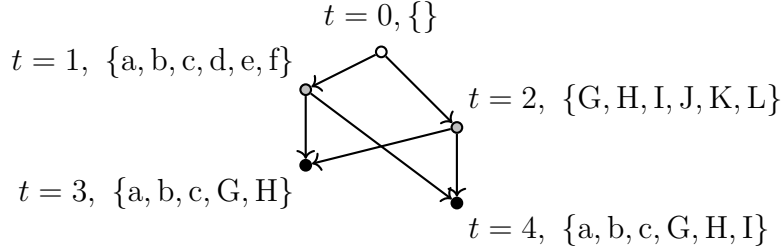
**Example 3.3.** Fig. 4 illustrates a history where two recombinants have the same parents. For the time function shown, the novelty profile is (5,1). The small second entry reflects the fact the two recombinant genomes are genetically close to one another. Exchanging the time values of the bottom-most two vertices yields another time function for this history, for





**Figure 3**

which the temporal novelty profile is (6,1). If we take the time values of the bottom-most two vertices to both be 3, then the temporal novelty profile is (6,5).



**Figure 4.** The temporal novelty profile of the history and time function shown is (5,1); the topological novelty profile is the same.

**Remark 3.4** (Stability). Suppose we are given histories  $\mathcal{E}$  and  $\mathcal{E}'$  indexed by the same phylogenetic graph  $G$  with  $d(\mathcal{E}_v, \mathcal{E}'_v) \leq \epsilon$  for all vertices  $v$  of  $G$ . Let  $t$  be any time function on the vertices of  $G$ . We then have that

$$d_\infty(\mathcal{N}(\mathcal{E}, t), \mathcal{N}(\mathcal{E}', t)) \leq \epsilon,$$

where for vectors  $A$  and  $B$  of the same length,  $d_\infty(A, B) := \max_i |A_i - B_i|$ . Thus, the temporal novelty profile is stable with respect to genomic perturbations.

However, by considering the history of Fig. 4, we can see that the temporal novelty profile is unstable with respect to perturbations of the time function.

### 3.2. The Topological Novelty Profile.

**Definition 3.5** (Relative Minimum Spanning Tree). Given a weighted graph  $G$  and a forest  $F \subset G$  (i.e., a vertex-disjoint collection of subtrees), we define a *spanning tree of  $G$  rel  $F$*  simply to be a spanning tree  $T$  of  $G$  containing  $F$ . We say  $T$  is a *minimum* spanning tree of  $G$  rel  $F$  if the sum of the edge weights of  $T$  is as small as possible, among all spanning trees of  $G$  rel  $F$ .

Note that by collapsing each tree in  $F$  to a point, the problem of finding a minimum spanning tree rel  $F$  is equivalent to the standard problem of finding an ordinary minimum spanning tree on a multigraph. (A multigraph is a graph which is allowed to have multiple edges between pairs of vertices.) Thus, all the standard facts about minimum spanning trees have analogues for relative minimum spanning trees. For example, we have the following:

**Proposition 3.6.** *A spanning tree  $T$  rel  $F$  is minimum if and only if for all  $i$ , the  $i^{\text{th}}$  largest edge weight is smaller than the  $i^{\text{th}}$  largest edge weight in any other spanning tree rel  $F$ .*

It follows from Proposition 3.6 that the collection of edge weights in a relative minimum spanning tree is independent of the choice of the tree.

**Definition 3.7** (Topological Novelty Profile). For  $\mathcal{E}$  a history indexed by a phylogenetic graph  $G$ , let  $F^G$  be the forest in  $G$  obtained by removing all edges pointing to recombinants. Let  $\bar{G}$  denote the complete graph with same vertex set as  $G$ . Regard  $\bar{G}$  as a weighted graph by taking the weight of edge  $(a, b)$  to be  $d(\mathcal{E}_a, \mathcal{E}_b)$ .

Let  $T$  be a minimum spanning tree of  $\bar{G}$  rel  $F^G$ . We define  $\mathcal{T}(\mathcal{E})$ , the *topological novelty profile* of  $\mathcal{E}$ , to be the list of distances

$$\{d(\mathcal{E}_a, \mathcal{E}_b) \mid (a, b) \in T \setminus F^G\},$$

counted with multiplicity and sorted in descending order.

We will observe in Section 5.2 that the topological novelty profile has an interpretation in terms of persistent homology.

Given two lists of numbers  $A$  and  $B$ , each sorted in decreasing order, we write  $A \leq B$  if  $|A| \leq |B|$  and for each  $i \in \{1, \dots, |A|\}$ ,  $A_i \leq B_i$ .

**Proposition 3.8.** *For any history  $\mathcal{E}$  with time function  $t$ ,*

$$\mathcal{T}(\mathcal{E}) \leq \mathcal{N}(\mathcal{E}, t).$$

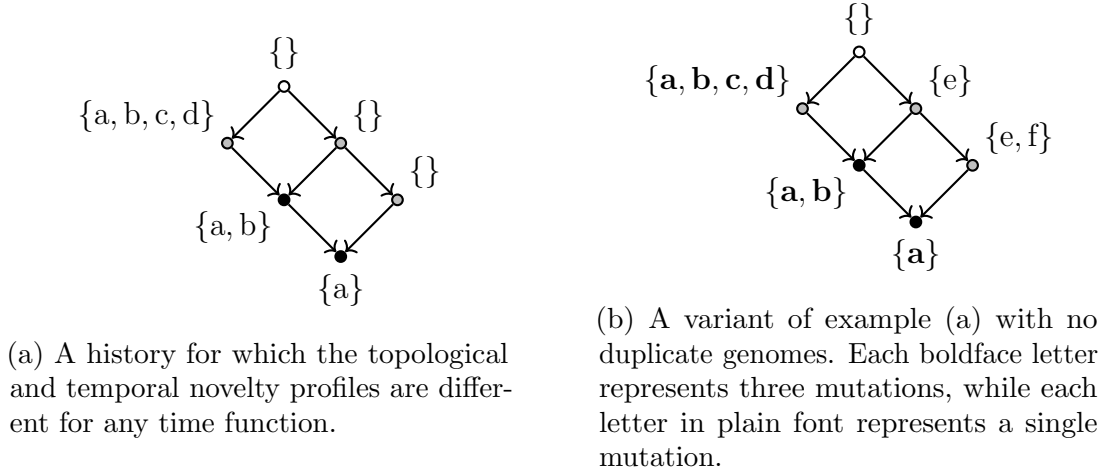
*That is, the topological novelty profile is a lower bound for the temporal novelty profile.*

*Proof.* Suppose  $\mathcal{E}$  is indexed by  $G$ . We construct a spanning tree  $T$  of  $\bar{G}$  rel  $F^G$  such the weights of edges in  $T \setminus F^G$  correspond to the temporal novelty profile. The result then follows from Proposition 3.6.

To construct  $T$ , for each recombinant  $r \in G$ , choose a vertex  $v(r)$  in  $G$  with  $t(v(r)) < t(r)$ , such that  $d(\mathcal{E}_{v(r)}, \mathcal{E}_r)$  is as small as possible among all such vertices. We take  $T$  to be the graph obtained from  $F^G$  by adding in the edge  $(v(r), r)$  for each recombinant  $r$ . It is easy to check that  $T$  is in fact a tree.  $\square$

**Example 3.9.** For the histories of Fig. 3a and Fig. 3b, the topological novelty profile is equal to the temporal one for all time functions. For the history and time function of Fig. 4, the topological and temporal novelty profiles are also equal, but one can select a different time function so that the two novelty profiles are not equal.

**Example 3.10.** Fig. 5a illustrates a history for which the temporal and topological novelty profiles are unequal for any choice of time function. The topological novelty profile is (1,1), whereas the temporal novelty profile is always (2,1). This example is degenerate, in the sense that the same genome (the empty one) appears at multiple vertices; Fig. 5b shows a variant of the example without this degeneracy.



**Figure 5**

Like temporal novelty profiles, topological novelty profiles are stable with respect to perturbations of the genomic data; we show this in Proposition 5.6. Proposition 4.5 below tells us that when  $G$  is a galled tree, the temporal and topological novelty profiles are in fact equal.

#### 4. HISTORIES INDEXED BY GALLED TREES

Our main results on the estimation of novelty profiles concern the special case that our phylogenetic graphs are galled trees.

The definition of galled tree we give is equivalent to the one given in [33, Definition 6.11.1]. As noted in [33], this is slightly more general than the original definition [29, 51], which requires the cycles in a galled tree to be node-disjoint.

**Definition 4.1** (Directed Loop). We say an undirected graph is a *loop* if its geometric realization is homeomorphic to a circle. We call a directed graph  $G$  a *directed loop* if

1. The undirected graph underlying  $G$  is a loop.
2.  $G$  has a unique source and unique sink.

**Definition 4.2** (Sum of Directed Graphs). For directed graphs  $G$  and  $H$ , with  $v$  a source in  $G$  and  $w$  any vertex in  $H$ , we define a directed graph  $G \vee_{v,w} H$  by taking the disjoint union of  $G$  and  $H$  and then identifying  $v$  and  $w$  (i.e., “gluing”  $v$  to  $w$ ). We call  $G \vee_{v,w} H$  a *sum* of  $G$  and  $H$ . (We do *not* define the sum  $G \vee_{v,w} H$  in the case that neither of the vertices  $v$  or  $w$  is a source.) We will sometimes write  $G \vee_{v,w} H$  simply as  $G \vee H$ , suppressing  $v, w$ .

**Definition 4.3** (Galled Tree). Let  $\mathcal{A}$  be the smallest collection of directed acyclic graphs such that:

1. Each rooted tree is in  $\mathcal{A}$ .
2. Each directed loop is in  $\mathcal{A}$ .
3. if  $G$  and  $H$  are in  $\mathcal{A}$ , then so is each sum  $G \vee H$ .

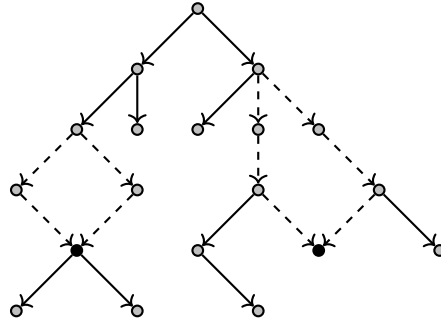
We define a *galled tree* to be a graph isomorphic to one in  $\mathcal{A}$ . Thus, informally, a galled tree is a graph obtained by iteratively gluing rooted trees and directed loops along single vertices, using the sum operation specified above.

We omit the easy proof of the following:

**Proposition 4.4.** *Any galled tree is a phylogenetic graph.*

Note that the recombinants in a galled tree  $G$  are in bijective correspondence with the directed loops in  $G$ .

Fig. 6 gives an example of a galled tree. It can be checked that the phylogenetic graph of Fig. 1 is not a galled tree.



**Figure 6.** A galled tree which can be constructed as the iterated sum of four rooted trees (solid edges) and two directed loops (dashed edges). The two recombinants are shown in black.

**Proposition 4.5** (Equality of Temporal and Topological Novelty Profiles on Galled Trees). *For any history  $\mathcal{E}$  indexed by a galled tree and time function  $t$ ,*

$$\mathcal{N}(\mathcal{E}, t) = \mathcal{T}(\mathcal{E}).$$

*Proof.* Suppose  $\mathcal{E}$  is indexed by the galled tree  $G$ . We use the notation from Definition 3.7. As in the case of ordinary minimum spanning trees, a minimum spanning tree of  $\bar{G}$  rel  $F$  can be constructed greedily, by considering the edges of  $\bar{G} \setminus F$  in order of increasing weight. In this construction, each edge of  $\bar{G} \setminus F$  added to the relative minimum spanning tree can be chosen to connect a recombinant  $r$  to a vertex  $v$  of the directed loop in  $G$  that has  $r$  as its sink. We then have that  $t(v) < t(r)$  and  $d(\mathcal{E}_v, \mathcal{E}_r) \leq d(\mathcal{E}_w, \mathcal{E}_r)$  for any other vertex  $w$  with  $t(w) < t(r)$ . Clearly, in this construction we never take  $r$  to be the same recombinant more than once. The result follows.  $\square$

**Remark 4.6** (Galled Trees as Models for Evolution in the Low-Recombination Limit). Given a probabilistic model generating a phylogenetic graph, one may ask what the probability is of obtaining a galled tree. This problem has previously been studied by simulation in [4], for a coalescent model of evolution. In Appendix A, we study the same problem analytically. We show that the problem reduces to the study of a finite-state Markov chain. A simple analysis of this Markov chain yields, for fixed population size  $n$ , a system of linear equations  $L(\rho)$  depending on a recombination rate parameter  $\rho$ , whose solution gives the probability  $P(n, \rho)$

of obtaining a galled tree. Solving these linear systems numerically for various values of  $\rho$  and  $n$ , we observe that as  $\rho$  tends to 0,  $P(n, \rho)$  tends to 1.

This indicates that histories indexed by galled trees are biologically reasonable models of evolution in low-recombination settings. While, from a biological standpoint, the specific bounds on  $\rho$  needed to obtain a galled tree with high probability are rather stringent in general, we do expect these bounds to hold in some settings of interest; see Section 9 for further discussion of this.

Regardless, from a mathematical perspective, the special case of galled trees seems to be a natural place to begin fleshing out theoretical foundations for the topological study of evolution.

## 5. TOPOLOGICAL PRELIMINARIES

In this section, we briefly review persistent homology and the related topological definitions and results we will need in the remainder of the paper. As a first application, we observe that the topological novelty profile admits a description in terms of persistent homology, and is therefore stable. We also briefly review some ideas from discrete Morse theory.

We assume that the reader is familiar with some standard concepts from elementary algebraic topology, including simplicial complexes, homology, and homotopy equivalence. Good introductions can be found in many places, e.g., [31, 40].

**5.1. Persistent Homology.** Our treatment of persistent homology will be terse; for a more thorough introduction to these ideas, including a discussion of some of the many applications of persistent homology to data analysis, see the surveys and textbooks [10, 11, 22, 43].

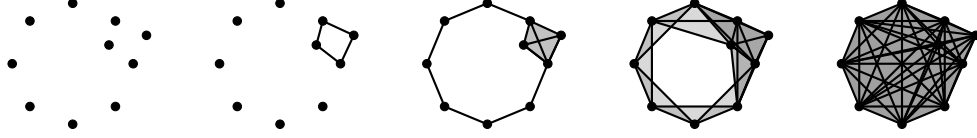
*Filtrations.* A *filtration* is a collection of topological spaces  $\{\mathcal{F}_r\}_{r \in [0, \infty)}$  such that  $\mathcal{F}_r \subset \mathcal{F}_s$  whenever  $r \leq s$ . A morphism  $f : \mathcal{F} \rightarrow \mathcal{G}$  of filtrations is a collection of continuous maps  $\{f_r : \mathcal{F}_r \rightarrow \mathcal{G}_r\}_{r \in [0, \infty)}$  such that the following diagram commutes for all  $r \leq s$ :

$$\begin{array}{ccc} \mathcal{F}_r & \hookrightarrow & \mathcal{F}_s \\ f_r \downarrow & & \downarrow f_s \\ \mathcal{G}_r & \hookrightarrow & \mathcal{G}_s \end{array}$$

We say  $f$  is an *objectwise homotopy equivalence* if each  $f_r$  is a homotopy equivalence. Intuitively, if two filtrations are connected by an objectwise homotopy equivalence, we should think of them as topologically equivalent; for further discussion of this point in the context of topological data analysis, see [7].

*Vietoris-Rips Filtrations.* For  $S$  a simplicial complex, we use square brackets to denote simplices of  $S$ . Thus, for example, the set of simplices of a triangle with vertex set  $\{a, b, c\}$  is  $\{[a], [b], [c], [ab], [bc], [ac]\}$ .

For  $P$  a finite metric space and  $r \in [0, \infty)$ , the *Vietoris-Rips complex* of  $P$  with scale parameter  $r$ , denoted  $\mathcal{V}(P)_r$ , is the simplicial complex with vertices  $P$  that contains simplex  $[p_1, p_2, \dots, p_n]$  if and only if  $\text{diameter}\{p_1, p_2, \dots, p_n\} \leq 2r$ . If  $r \leq s$ , then  $\mathcal{V}(P)_r \subset \mathcal{V}(P)_s$ , so  $\mathcal{V}(P) := \{\mathcal{V}(P)_r\}_{r \in [0, \infty)}$  is a filtration; see Fig. 7.



**Figure 7.** Rips complexes  $\mathcal{V}(P)_r$  on a point simple point cloud  $P \subset \mathbb{R}^2$ , for several choices of  $r$ .

*Persistence Modules.* A *persistence module*  $M$  consists of a collection of vector spaces  $\{M_r\}_{r \in [0, \infty)}$ , together with a collection of linear maps  $\{M_{r,s} : M_r \rightarrow M_s\}_{r \leq s}$  such that

1. for all  $r \leq s \leq t$  the following diagram commutes:

$$\begin{array}{ccc} M_r & & \\ M_{r,s} \downarrow & \searrow M_{r,t} & \\ M_s & \xrightarrow{M_{s,t}} & M_t. \end{array}$$

2.  $M_{r,r} = \text{id}_{M_r}$  for all  $r$ .

We say  $M$  is *pointwise finite dimensional (p.f.d.)* if  $\dim M_r < \infty$  for all  $r$ .

Similar to the definition for filtrations, a morphism  $f : M \rightarrow N$  of persistence modules is a collection of linear maps  $\{f_r : M_r \rightarrow N_r\}_{r \in [0, \infty)}$  such that for all  $r \leq s$ , the following diagram commutes:

$$\begin{array}{ccc} M_r & \xrightarrow{M_{r,s}} & M_s \\ f_r \downarrow & & \downarrow f_s \\ N_r & \xrightarrow{N_{r,s}} & N_s. \end{array}$$

We say  $f$  is an isomorphism if each of the maps  $f_r$  is an isomorphism.

*Direct Sums of Persistence Modules.* We assume that the reader is familiar with the definition of the direct sum of vector spaces from linear algebra. For linear maps  $f : V_1 \rightarrow W_1$  and  $g : V_2 \rightarrow W_2$ , we define the direct sum

$$f \oplus g : V_1 \oplus V_2 \rightarrow W_1 \oplus W_2$$

by taking  $f \oplus g(v, w) = (f(v), g(w))$ . We then define the sum  $M \oplus N$  to be the persistent module given by

$$(M \oplus N)_r = M_r \oplus N_r, \quad (M \oplus N)_{r,s} = M_{r,s} \oplus N_{r,s}.$$

More generally, we can define the direct sum of an arbitrary collection of persistence modules in the same way.

*Reduced Homology.* Fix a field  $K$ . (For example, we can take  $K = \mathbb{Q}$ , or  $K = \mathbb{Z}_2$ , the field with two elements.) For  $i \geq 0$ , let  $H_i$  denote the  $i^{\text{th}}$  reduced singular homology functor with coefficients in  $K$ . Thus,  $H_i$  maps each topological space  $S$  to a  $K$ -vector space  $H_i(S)$ , and maps each continuous function  $f : S \rightarrow T$  to a linear map  $f_* : H_i(S) \rightarrow H_i(T)$ . Applying  $H_i$

to each space and each inclusion map in a filtration  $\mathcal{F}$  gives us a persistence module  $H_i(\mathcal{F})$ . Moreover, a morphism of filtrations  $f : \mathcal{F} \rightarrow \mathcal{G}$  induces a morphism  $f_* : H_i(\mathcal{F}) \rightarrow H_i(\mathcal{G})$ .

**Lemma 5.1.** *If a morphism of filtrations  $f : \mathcal{F} \rightarrow \mathcal{G}$  is an objectwise homotopy equivalence, then for any  $i \geq 0$ ,  $f_* : H_i(\mathcal{F}) \rightarrow H_i(\mathcal{G})$  is an isomorphism.*

*Proof.* It is a standard fact that if a continuous map  $g$  is a homotopy equivalence, then  $H_i(g)$  is an isomorphism. This gives the result.  $\square$

*Barcodes.* We say  $\mathcal{J} \subset \mathbb{R}$  is an *interval* if  $\mathcal{J}$  is nonempty and connected. For  $\mathcal{J}$  an interval, define the *interval module*  $I^{\mathcal{J}}$  to be the persistence module such that

$$I_r^{\mathcal{J}} = \begin{cases} K & \text{if } r \in \mathcal{J}, \\ 0 & \text{otherwise.} \end{cases} \quad I_{r,s}^{\mathcal{J}} = \begin{cases} \text{id}_K & \text{if } r \leq s \in \mathcal{J}, \\ 0 & \text{otherwise.} \end{cases}$$

**Theorem 5.2** (Structure of Persistence Modules [19]). *If  $M$  is a p.f.d. persistence module, then there exists a unique collection of intervals  $\mathcal{B}_M$  such that*

$$M \cong \bigoplus_{\mathcal{J} \in \mathcal{B}_M} I^{\mathcal{J}}.$$

We call  $\mathcal{B}_M$  the *barcode* of  $M$ . For  $\mathcal{F}$  a filtration, we write  $\mathcal{B}_{H_i(\mathcal{F})}$  simply as  $\mathcal{B}_i(\mathcal{F})$ . Similarly, for  $P$  a finite metric space, we write  $\mathcal{B}_i(\mathcal{V}(P))$  simply as  $\mathcal{B}_i(P)$ .

**Definition 5.3** (Gromov-Hausdorff Distance). Given two subspaces  $P, Q$  of a metric space  $Z$ , we define the *Hausdorff distance* between  $P$  and  $Q$ , by

$$d_H(P, Q) := \max\left\{\sup_{p \in P} \inf_{q \in Q} d(p, q), \sup_{q \in Q} \inf_{p \in P} d(p, q)\right\}$$

For  $P$  and  $Q$  any compact metric spaces, define  $d_{GH}(P, Q)$ , the *Gromov-Hausdorff distance* between  $P$  and  $Q$ , to be the infimum of  $d_H(\gamma(P), \kappa(Q))$  over all isometric embeddings  $\gamma : P \rightarrow Z$ ,  $\kappa : Q \rightarrow Z$  into a third metric space  $Z$ .

The following stability result is well known, and plays a central result in topological data analysis. Let  $d_B$  denote the *bottleneck distance* on persistence barcodes, as defined for example in [16].

**Theorem 5.4** (Stability of Persistent Homology [14, 16, 17]). *For any finite metric spaces  $P, Q$  and  $i \geq 0$ ,*

$$d_B(\mathcal{B}_i(P), \mathcal{B}_i(Q)) \leq d_{GH}(P, Q).$$

The following variant of Theorem 5.4, appearing in [30], can be proven by a slight modification of the proof of Theorem 5.4.

**Theorem 5.5** (Stability for a Metric Subspace [30]). *For finite metric spaces  $P \subset Q$  and  $i \geq 0$ ,*

$$d_B(\mathcal{B}_i(P), \mathcal{B}_i^S(Q)) \leq \frac{1}{2} d_H(P, Q),$$

where  $\mathcal{B}_i^S(Q)$  is the barcode obtained by shifting each interval of  $\mathcal{B}_i(Q)$  to the right by  $\frac{1}{2}d_H(P, Q)$ .

**5.2. Stability of the Topological Novelty Profile.** For  $\mathcal{E}$  a history indexed by  $G$  and  $F^G \subset G$  the forest of Definition 3.7, define a filtration  $\mathcal{F}$  by  $\mathcal{F}_r := F^G \cup \mathcal{V}(M)_r$ . It's easy to check that  $\mathcal{T}(\mathcal{E})$ , the topological novelty profile of  $\mathcal{E}$ , is exactly the list of right endpoints of intervals in  $\mathcal{B}_0(\mathcal{F})$ , possibly with some copies of 0 added in.

We can use this description to obtain a simple stability result for topological novelty profiles analogous to the one for temporal novelty profiles mentioned in Remark 3.4:

**Proposition 5.6.** *Given histories  $\mathcal{E}$  and  $\mathcal{E}'$  indexed by the same phylogenetic graph  $G$  with  $d(\mathcal{E}_v, \mathcal{E}'_v) \leq \epsilon$  for all vertices  $v$  of  $G$ , we have*

$$d_\infty(\mathcal{T}(\mathcal{E}), \mathcal{T}(\mathcal{E}')) \leq \epsilon.$$

*Proof.* This follows immediately from a generalized version of the stability theorem for persistent homology, as described in [5, 13, 15].  $\square$

**5.3. Discrete Morse Theory.** The proof of our main results relies on discrete Morse theory (DMT), a well known combinatorial theory concerning topology-preserving collapses of cell complexes. We will not need the full strength of standard DMT; we review only what we need for our proof. See [26] for a detailed introduction to DMT.

For  $S$  a simplicial complex, the *Hasse graph*  $G_S$  of  $S$  is the directed graph with vertices the simplices of  $S$  and an edge from  $s$  to  $s'$  if and only if  $s'$  is a codimension-1 face of  $s$ . We say that a matching  $X$  in  $G_S$  is *acyclic* if when we modify the graph  $G_S$  by reversing the orientation of all edges in  $X$ , while leaving the orientation of all other edges unchanged, we obtain an directed acyclic graph.

A *discrete gradient vector field (DGVF)*  $X$  on  $S$  is an acyclic matching in  $G_S$ . We say a simplex  $\sigma \in S$  is *critical* in  $X$  if  $\sigma$  is not matched in  $X$ .

The acyclicity condition admits an alternative formulation which is often convenient. Given a matching  $X$  in  $G_S$ , we define an  $X$ -*path* to be a sequence of simplices in  $S$

$$\sigma_0, \tau_0, \sigma_1, \tau_1, \dots, \sigma_m, \tau_m, \sigma_{m+1}$$

such that for each  $j \in \{0, \dots, m\}$ , the following are true:

- $\sigma_j$  is a face of  $\tau_j$  and  $X$  matches  $\sigma_j$  to  $\tau_j$ ,
- $\sigma_{j+1}$  is a codimension-1 face of  $\tau_j$ ,
- $\sigma_j \neq \sigma_{j+1}$ .

We say the  $X$ -path is a *non-trivial* if  $m \geq 0$ , and *closed* if  $\sigma_0 = \sigma_{m+1}$ .

**Proposition 5.7** ([26, Theorem 6.2]). *A matching  $X$  in  $G_S$  is a DGVF if and only if there exists no non-trivial closed  $X$ -path.*

The following is one of the basic results of discrete Morse theory:

**Proposition 5.8** ([25, 37]).

- (i) *Suppose that  $X$  is a DGVF on a finite simplicial complex  $S$ . Then  $S$  is homotopy equivalent to a CW-complex with exactly one cell of dimension  $i$  for each critical  $i$ -simplex of  $X$ .*
- (ii) *If the critical simplices of  $X$  form a subcomplex  $S' \subset S$ , then in fact  $S$  deformation retracts onto  $S'$ .*



## 6. BARCODES OF HISTORIES INDEXED BY GALLED TREES

The topological novelty profile and 0<sup>th</sup> persistence barcode of an evolutionary history are closely related by the following result, whose easy verification we leave to the reader:

**Proposition 6.1.** *Suppose we are given a history  $\mathcal{E}$  and  $\delta > 0$  such that  $d(\mathcal{E}_v, \mathcal{E}_w) < \delta$  whenever  $w$  is a clone with parent  $v$ . Then the lists obtained from  $\text{lengths}(\mathcal{B}_0(\mathcal{E}))$  and  $\mathcal{T}(\mathcal{E})$  by removing all entries less than  $\delta$  are equal.*

This suggests that in some cases, 0<sup>th</sup> barcodes may be useful in the study of recombination. However, in cases where the minimum  $\delta$  satisfying the condition of Proposition 6.1 is large, or where we only have a subsample of the history, 0<sup>th</sup> barcodes may not offer useful information. This, together with the earlier theoretical result of Chan, Rabadan and Carlsson relating recombination to higher persistence barcodes (Theorem 6.4 below), motivates us to consider the relationship between the topological novelty profile and the higher barcodes of a history.

In this section, we present our main result relating barcodes to novelty profiles in the galled tree setting (Theorem 6.18). The technical heart of our proof is a topological description of the Vietoris-Rips filtration of an *almost linear* metric space, which we give in Section 6.3. Our arguments rely heavily on discrete Morse theory.

**6.1. Barcodes of Histories indexed by Trees.** We first review the key result of Chan et al. on the barcodes of histories indexed by trees.

**Definition 6.2** (Tree-Like Metric Space). We call an undirected tree with a non-negative weight function on its edges a *weighted tree*. A metric space  $P$  is called *tree-like* if it is isometric to a subspace of a metric space arising from the shortest-path metric on a weighted tree.

**Proposition 6.3** ([12]). *If  $P$  is a tree-like metric space, then for all  $r \in [0, \infty)$ , each component of  $\mathcal{V}(P)_r$  is contractible. Hence,  $\mathcal{B}_i(P) = 0$  for  $i \geq 1$ .*

In [12], only the part of Proposition 6.3 about triviality of barcodes is stated, and not the stronger contractibility result. However, the contractibility result follows immediately from the proof given in [12], using the nerve theorem [31] in place of a Mayer-Vietoris argument.

The following result, due to Chan et al., makes precise the idea that for  $i \geq 1$ , a non-empty barcode  $\mathcal{B}_i(S)$  serves as a certificate that recombination is present in the history from which  $S$  was sampled.

**Theorem 6.4** ([12]). *If  $G$  is a tree,  $\mathcal{E}$  is a history indexed by  $G$ , and  $S \subseteq \mathcal{E}$ , then  $\mathcal{B}_i(S) = \emptyset$  for  $i \geq 1$ .*

*Proof.* If  $S$  is a subset of a history indexed by a tree, then  $\text{met } S$  is easily seen to be tree-like. Hence, the result follows from Proposition 6.3.  $\square$

**Remark 6.5.** In the absence of recombination, homoplasies (recurrent mutations that violate the infinite sites assumption) can lead to a metric space that is not tree-like. However, as indicated in Remark 2.7, a small number of homoplasies causes a correspondingly small deviation from tree-likeness (with respect to Gromov-Hausdorff distance). A single recombination event, on the other hand, can yield a metric space that is arbitrarily far from a tree-like one.

**6.2. Metric Decomposition of an Evolutionary History.** Define a *based metric space* simply to be a metric space  $P$ , together with a choice of basepoint  $p \in P$ .

**Definition 6.6** (Sum of Based Metric Spaces). For based metric spaces  $P$  and  $Q$  with basepoints  $p \in P$ ,  $q \in Q$ , we regard the wedge sum  $P \vee Q$  as a metric space, with the metric given by

$$d_{P \vee Q}(x, y) = \begin{cases} d_P(x, y) & \text{if } x, y \in P, \\ d_Q(x, y) & \text{if } x, y \in Q, \\ d_P(x, p) + d_Q(q, y) & \text{if } x \in P, y \in Q. \end{cases}$$

For based metric spaces  $P$  and  $Q$ , let  $\mathcal{V}(P) \vee \mathcal{V}(Q)$  denote the *wedge sum filtration*, given by

$$(\mathcal{V}(P) \vee \mathcal{V}(Q))_r := \mathcal{V}(P)_r \vee \mathcal{V}(Q)_r.$$

**Proposition 6.7.** *For finite based metric spaces  $P$  and  $Q$ , the inclusion*

$$\mathcal{V}(P) \vee \mathcal{V}(Q) \hookrightarrow \mathcal{V}(P \vee Q)$$

*is an objectwise homotopy equivalence. In particular, for any  $i \geq 0$ ,*

$$\mathcal{B}_i(P \vee Q) = \mathcal{B}_i(P) \cup \mathcal{B}_i(Q).$$

*Proof.* We give a proof using discrete Morse theory. For  $r \in \mathbb{R}$ , define a DGVF on  $\mathcal{V}(P \vee Q)_r$  as follows: For a simplex  $\sigma$  in  $\mathcal{V}(P \vee Q)_r$  containing vertices in both  $P$  and  $Q$ , but not the common vertex  $p = q$ , we match  $\sigma$  to  $\{p = q\} \cup \sigma$ . It is clear that this matching is acyclic, hence indeed gives a well-defined DGVF whose set of critical simplices is  $(\mathcal{V}(P) \vee \mathcal{V}(Q))_r$ . Thus, by Proposition 5.8 (ii), the inclusion

$$(\mathcal{V}(P) \vee \mathcal{V}(Q))_r \rightarrow \mathcal{V}(P \vee Q)_r$$

is a homotopy equivalence.

We now check that  $\mathcal{B}_i(P \vee Q) = \mathcal{B}_i(P) \cup \mathcal{B}_i(Q)$ . By Lemma 5.1,  $\mathcal{B}_i(P \vee Q) = \mathcal{B}_i(\mathcal{V}(P) \vee \mathcal{V}(Q))$ . Moreover,  $\mathcal{B}_i(\mathcal{V}(P) \vee \mathcal{V}(Q)) = \mathcal{B}_i(P) \cup \mathcal{B}_i(Q)$ ; this follows from the corresponding result for topological spaces [31, Corollary 2.25], using the fact that the isomorphism appearing in the statement of that result is natural.  $\square$

**Remark 6.8.** Proposition 6.7 has a category theoretic-interpretation: It says that reduced persistent homology commutes with coproducts in the categories of based metric spaces and persistence modules, where morphisms of metric spaces are 1-Lipschitz maps sending basepoint to basepoint.

**Remark 6.9.** Proposition 6.7 has also been discovered independently by the authors of [2]. Their work also establishes the result for infinite based metric spaces and for Čech filtrations.

We leave the easy verification of the following to the reader:

**Proposition 6.10.** *Suppose  $G$  is a phylogenetic graph with  $G = G^1 \vee G^2$  for subgraphs  $G^1, G^2 \subseteq G$ ,  $\mathcal{E}$  is a history indexed by  $G$ , and  $\mathcal{E}^1$  and  $\mathcal{E}^2$  are the respective restrictions of  $\mathcal{E}$  to  $G^1$  and  $G^2$ . Then*

$$\text{met } \mathcal{E} \cong \text{met } \mathcal{E}^1 \vee \text{met } \mathcal{E}^2.$$

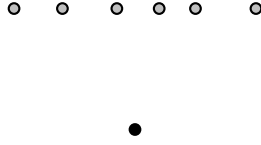
**Theorem 6.11.** Suppose a galled tree  $G$  is an iterated sum of directed loops  $G^1, \dots, G^k \subset G$  and rooted trees  $G^{k+1}, \dots, G^l \subset G$ , and that  $\mathcal{E}$  is a history indexed by  $G$ . Let  $\mathcal{E}^j$  denote the restriction of  $\mathcal{E}$  to  $G^j$ .

- (i) There is an objectwise homotopy equivalence from an iterated wedge sum of the filtrations  $\mathcal{V}(E^j)$  to  $\mathcal{V}(\mathcal{E})$ .
- (ii) For  $i \geq 1$ ,

$$\mathcal{B}_i(\mathcal{E}) = \bigcup_{j=1}^k \mathcal{B}_i(\mathcal{E}^j).$$

*Proof.* (i) follows from Propositions 6.7 and 6.10. (ii) follows from (i) and Proposition 6.3.  $\square$

**6.3. Vietoris-Rips Filtrations of Almost Linear Metric Spaces.** As mentioned in the introduction, we say a non-empty finite metric space  $P$  is *almost linear* if there is a point  $p \in P$  such that  $P \setminus \{p\}$  is isometric to a finite subset of  $\mathbb{R}$ . We call any such point  $p$  a *distinguished point*. See Fig. 8.



**Figure 8.** An almost linear metric space embedded in  $\mathbb{R}^2$ . The unique distinguished point is shown in solid black. Note that not all almost linear metric spaces can be embedded in  $\mathbb{R}^2$ .

**Proposition 6.12.** If  $\mathcal{E}$  is a history indexed by a directed loop, then  $\text{met } \mathcal{E}$  is almost linear.

*Proof.* Let  $p$  be the unique recombinant.  $\text{met } \mathcal{E} \setminus \{\mathcal{E}_p\}$  is isometric to a subset of  $\mathbb{R}$ .  $\square$

In view of Theorem 6.11 and Proposition 6.12, to understand the topology of Vietoris-Rips filtrations of histories indexed by galled trees, it suffices to understand the topology of Vietoris-Rips filtrations of almost linear metric spaces. We now describe the latter:

**Theorem 6.13** (Topology of the Vietoris-Rips Filtration of an Almost Linear Metric Space). Let  $P$  be an almost linear metric space with distinguished point  $p$ .

- (i) For each  $r \in [0, \infty)$ , the connected component  $C_r$  of  $\mathcal{V}(P)_r$  containing  $p$  is either contractible or homotopy equivalent to a circle, and each other component of  $\mathcal{V}(P)_r$  is contractible. In particular,  $\mathcal{B}_i(P) = 0$  for  $i \geq 2$ .
- (ii) If  $C_r$  and  $C_{r'}$  are both homotopy equivalent to circles and  $r \leq r'$ , then the inclusion  $C_r \hookrightarrow C_{r'}$  is a homotopy equivalence. Thus,  $\mathcal{B}_1(P)$  has at most one interval.
- (iii) The unique interval of  $\mathcal{B}_1(P)$ , when it exists, has length at most  $d(p, P \setminus \{p\})$  and is contained in the interval

$$[d(p, P \setminus \{p\}), \text{diameter}(P \setminus \{p\})/2).$$

**Remark 6.14.** Together, Theorem 6.11 (i), Proposition 6.12, and Theorem 6.13 (i) tell us that for  $\mathcal{E}$  a history indexed by a galled tree  $G$  and  $r \in [0, \infty)$ , each component of  $\mathcal{V}(\mathcal{E})_r$  is homotopy equivalent to a bouquet of circles.

**Remark 6.15.** In analogy with the definition of an almost linear metric space, we can define an *almost tree-like* metric space to be one obtained from a tree-like metric space by adding a single point. We conjecture that Theorem 6.13 also holds for almost tree-like metric spaces.

As a first step towards proving Theorem 6.13, we use discrete Morse theory to show the following:

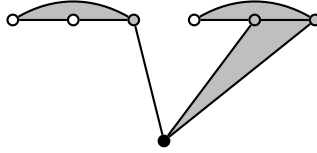
**Lemma 6.16.** *For any  $r \in [0, \infty)$ , each component of  $\mathcal{V}(P)_r$  is contractible or deformation retracts onto a wedge sum of finitely many circles.*

*Proof.* Any component of  $\mathcal{V}(P)_r$  not containing  $p$  is tree-like, and so is contractible by Proposition 6.3. Thus, we henceforth assume without loss of generality that  $\mathcal{V}(P)_r$  is connected.

By choosing an isometric embedding  $P \setminus \{p\} \hookrightarrow \mathbb{R}$ , we may regard  $P \setminus \{p\}$  as a subset of  $\mathbb{R}$ . Let  $P_{\text{left}} \subset P \setminus \{p\}$  denote the set of points  $y$  such that

1.  $y$  is not incident to  $p$ ,
2. there is no vertex of  $P \setminus \{p\}$  to the left of  $y$  lying in the same connected component of  $\mathcal{V}(P \setminus \{p\})_r$  and incident to  $p$ .

See Fig. 9 for an illustration of  $P_{\text{left}}$ .



**Figure 9.** Illustration of the Vietoris-Rips complex  $\mathcal{V}(P)_r$  for an almost linear metric space  $P$ , and some choice of scale parameter  $r$ . Here, the metric on  $P$  is not assumed to be the one given by the shown embedding of the points in the plane. The exceptional point is solid black, points of  $P_{\text{left}}$  are white, and the remaining points are gray.

We show via a simple discrete Morse theory argument that  $\mathcal{V}(P)_r$  deformation retracts onto  $\mathcal{V}(P \setminus P_{\text{left}})_r$ . Define a DGVF  $X$  on  $\mathcal{V}(P)_r$  as follows: For  $j \geq 2$  and

$$\sigma := [a_1 < a_2 < \dots < a_j]$$

a simplex in  $\mathcal{V}(P)_r$  such that  $a_1 \in P_{\text{left}}$  and  $a_2$  is the point in  $P$  immediately to the right of  $a_1$ ,  $X$  matches  $\sigma$  to its face  $[a_1, a_3, \dots, a_j]$ . To see that  $X$  is acyclic, note that for any  $X$ -path

$$\sigma_0, \tau_0, \dots, \sigma_m, \tau_m, \sigma_{m+1},$$

the  $\tau_j$  are strictly increasing with respect to the lexicographical order induced by the vertex ordering. If  $m \geq 0$  and  $\sigma_0 = \sigma_{m+1}$ , then

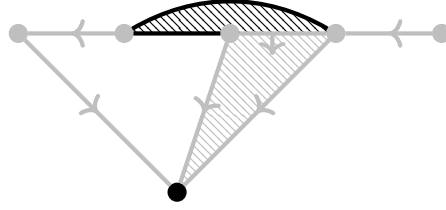
$$\sigma_0, \tau_0, \dots, \sigma_m, \tau_m, \sigma_0, \tau_0, \sigma_1$$

is an  $X$ -path with  $\tau_m < \tau_0$ , so there cannot exist a non-trivial closed  $X$ -path. Therefore  $X$  is acyclic.

Furthermore,  $X$  matches every simplex containing a point in  $P_{\text{left}}$ , so the critical simplices of  $X$  form the subcomplex  $\mathcal{V}(P \setminus P_{\text{left}})_r$ . Hence,  $\mathcal{V}(P)_r$  deformation retracts onto  $\mathcal{V}(P \setminus P_{\text{left}})_r$  by Proposition 5.8 (ii). We thus may assume without loss of generality that  $P_{\text{left}} = \emptyset$ .

We will define a DGVF  $Y$  on  $\mathcal{V} := \mathcal{V}(P)_r$  with a single critical 0-simplex, and no critical simplices of dimension greater than one; the result then follows from Proposition 5.8 (i). We define  $Y$  in two steps, first giving a simple definition of a DGVF  $X$  on  $\mathcal{V}$  and then refining this by matching more simplices. To start, we order the vertices in  $P$  by taking  $\{p\}$  to be the minimum, and ordering  $P - \{p\}$  from left to right, via the chosen embedding of  $P - \{p\}$  into  $\mathbb{R}$ . Henceforth, it will be our convention that the vertices of a simplex in  $\mathcal{V}$  are always written in increasing order.

We define  $X$  using a construction of Matt Kahle [35] (which in fact gives a DGVF on any simplicial complex with ordered vertex set): If a simplex  $\sigma = [a_1, a_2, \dots, a_j]$  has a coface  $a_0 \cup \sigma := [a_0, a_1, a_2, \dots, a_j]$  with  $a_0 < a_1$ ,  $X$  matches  $\sigma$  to  $a_0 \cup \sigma$  with  $a_0$  as small as possible.  $X$  matches no other simplices. It is easy to check that this in fact gives a well-defined DGVF. See Fig. 10 for an illustration.



**Figure 10.** Illustration of the discrete gradient vector field  $X$  on a Vietoris-Rips complex of an almost-linear metric space, with the bottom vertex ordered first, and the remaining vertices ordered left-to-right. Matched simplices are gray, and matched pairs are denoted with an arrow pointing away from the simplex of lower dimension. Critical simplices are black. Thus,  $X$  has a one critical 0-simplex, two critical 1-simplices, and one critical 2-simplex.

Clearly,  $[p]$  is critical in  $X$ , and since we assume  $P_{\text{left}} = \emptyset$ , no other vertex is critical. The following describes the remaining critical simplices in  $X$ :

**Lemma 6.17.** *For  $j \geq 2$ , a simplex  $[a_1, \dots, a_j]$  is critical in  $X$  if and only if the following three conditions are satisfied:*

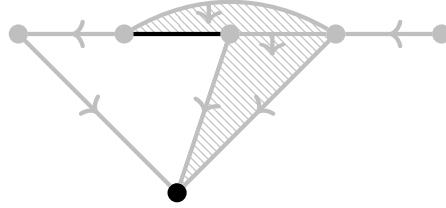
1.  $a_1 \neq p$ ,
2.  $[q, a_1, \dots, a_j] \notin \mathcal{V}$  for any  $q < a_1$ ,
3.  $[p, a_2, a_3, \dots, a_j] \in \mathcal{V}$ .

*Proof.* If all three conditions hold, then by condition 2,  $\sigma := [a_1, \dots, a_j]$  cannot be the simplex of lower dimension in a pair matched by  $X$ , and by condition 3,  $X$  matches  $[a_2, a_3, \dots, a_j]$  to  $[p, a_2, a_3, \dots, a_j]$ , so by condition 1,  $\sigma$  cannot be the simplex of higher dimension in a pair matched by  $X$ . Thus  $\sigma$  is critical in  $X$ .

Conversely, if  $\sigma$  is critical in  $X$ , then condition 1 holds, for else  $\sigma$  would match to  $[a_2, a_3, \dots, a_j]$ . Condition 2 holds, for else  $\sigma$  would match to a simplex of higher dimension. Finally, condition 3 holds, for else  $[a_2, a_3, \dots, a_j]$  would match to some simplex  $[q, a_2, \dots, a_j] \in \mathcal{V}$  with  $p < q < a_1$ , implying that  $[q, a_1, \dots, a_j] \in \mathcal{V}$ , and hence contracting the criticality of  $\sigma$ .  $\square$

Note that Lemma 6.17 implies in particular that if  $[a_1, \dots, a_j]$  is critical, then  $a_1$  is not incident to  $p$ , since otherwise, in view of condition 3, condition 2 would be violated.

The lemma suggests a way to extend  $X$  to a DGVF  $Y$  with the desired properties: For  $[a_1, a_2, a_3, \dots, a_j]$  a critical simplex in  $X$  with  $j \geq 3$ , suppose there exists no vertex  $b$  such that  $[p, b] \in \mathcal{V}$  and  $a_1 < b < a_2$ . It follows easily from Lemma 6.17 that  $[a_1, a_3, \dots, a_j]$  is also critical in  $X$ . We match  $[a_1, a_2, a_3, \dots, a_j]$  to  $[a_1, a_3, \dots, a_j]$  in  $Y$ . We take all matched pairs in  $Y \setminus X$  to be of this form. Fig. 11 illustrates the extension of the DGVF  $X$  of example Fig. 10 to the DGVF  $Y$ .



**Figure 11.** The extension of the DGVF  $X$  of example Fig. 10 to the DGVF  $Y$ .  $Y$  contains one pair of matched simplices not in  $X$ : The curved 1-simplex in the top of the figure now is matched with its coface. Thus,  $Y$  has two critical simplices: A critical 0-simplex and a critical 1-simplex.

It's easy to see that the resulting matching  $Y$  matches every critical simplex of  $X$ , except  $[p]$  and those 1-simplices of the form

$$\{[a, c] \mid \text{there does not exist } b \text{ with } a < b < c \text{ and } [p, b] \in \mathcal{V}\}.$$

In particular,  $Y$  has no unmatched simplices of dimension greater than one.

It remains only to check that  $Y$  is acyclic. We claim that in any  $Y$ -path

$$\sigma_0, \tau_0, \dots, \sigma_m, \tau_m, \sigma_{m+1},$$

no two distinct  $\tau_j$  are equal. From this, it follows that there does not exist a non-trivial closed  $Y$ -path, so  $Y$  is indeed acyclic. To verify the claim, we make three simple observations: Letting  $\tau_j^1$  denote the minimum vertex in  $\tau_j$ , we have that for any  $j \in \{0, \dots, m-1\}$ ,

1. If  $\tau_{j+1}$  is matched by  $X$ , then  $\tau_j^1 > \tau_{j+1}^1$ .
2. If  $\tau_{j+1}$  is matched by  $Y \setminus X$ , then  $\tau_j$  is matched by  $X$  and  $\tau_{j+1}^1 = \tau_j^1$ .
3.  $\tau_j \neq \tau_{j+1}$ .

By observations 1 and 2, we have that  $\tau_j^1 > \tau_k^1$  for all  $k \in \{j+2, j+3, \dots, m\}$ . The claim follows from this and observation 3.  $\square$

*Proof of Theorem 6.13 (i).* As in the proof of Lemma 6.16, we assume without loss of generality that  $\mathcal{V} := \mathcal{V}(P)_r$  is connected, and that  $P_{\text{left}} = \emptyset$ . The fundamental group of a wedge

sum of circles is free [31, Example 1.21], so by Lemma 6.16,  $\pi_1(\mathcal{V}, p)$  is free. To establish Theorem 6.13 (i), it suffices to show that  $\pi_1(\mathcal{V}, p)$  is trivial or cyclic.

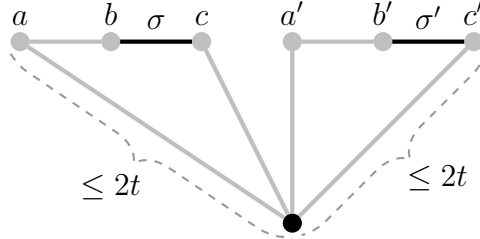
To show this, we first note that the argument of Lemma 6.16 yields a basis for  $\pi_1(\mathcal{V}, p)$ , as follows: Let  $\Gamma$  denote the set of critical 1-simplices of the DGVF  $Y$  defined in the proof of Lemma 6.16 above. For  $\sigma = [b, c] \in \Gamma$  with  $b < c$ , let  $a \in P - \{p\}$  denote the maximum vertex such that  $a < b$  and  $[p, a] \in \mathcal{V}$ . Such  $a$  always exists by our assumption that  $P_{\text{left}} = \emptyset$ . Let us regard  $S^1$  as a based topological space, with the basepoint denoted as 1, and let  $\gamma_\sigma : S^1 \rightarrow [p, a] \cup [a, c] \cup [p, c]$  be a homeomorphism sending 1 to  $p$ .

We now observe that  $G := \{\gamma_\sigma \mid \sigma \in \Gamma\}$  is a basis for  $\pi_1(\mathcal{V}, p)$ . For  $\sigma \in \Gamma$ , let  $S_\sigma^1$  denote a copy of  $S^1$ . The proof of Proposition 5.8 (i) presented in [25] gives a (not necessarily unique) homotopy equivalence  $h : \mathcal{V} \rightarrow \bigvee_{\sigma \in \Gamma} S_\sigma^1$  mapping the interior of  $\sigma$  homeomorphically to  $S_\sigma^1 \setminus \{1\}$ , so that  $h \circ \gamma_\sigma$  is homotopic either to the inclusion  $\iota_\sigma : S_\sigma^1 \hookrightarrow \bigvee_{\sigma \in \Gamma} S_\sigma^1$ , or to its inverse in  $\pi_1(\bigvee_{\sigma \in \Gamma} S_\sigma^1, 1)$ . Since  $h$  is a homotopy equivalence and  $\{\iota_\sigma \mid \sigma \in \Gamma\}$  is a basis for  $\pi_1(\bigvee_{\sigma \in \Gamma} S_\sigma^1, 1)$ , we see that  $G$  is a basis for  $\pi_1(\mathcal{V}, p)$ , as desired.

To finish the proof of Theorem 6.13 (i), it remains to show that  $|G| \leq 1$ . To do so, we apply the triangle inequality. Our argument is illustrated in Fig. 12. For  $[b, c] = \sigma \in \Gamma$  with  $b < c$ , let  $a < b$  be as above, and for  $[b', c'] = \sigma' \in \Gamma$  with  $b' < c'$ , define  $a' < b'$  in the same way. To arrive at a contradiction, suppose  $\sigma \neq \sigma'$ . Then either  $c \leq a'$  or  $c' \leq a$ . Switching the labels of  $\sigma$  and  $\sigma'$  if necessary, we may assume without loss of generality that  $c \leq a'$ . We have  $[p, a], [p, c'] \in \mathcal{V}$ , so  $d(a, p) \leq 2r$  and  $d(p, c') \leq 2r$ . By the triangle inequality,  $d(a, c') \leq 4r$ . Thus, since  $P \setminus \{p\}$  is isometric to a subset of  $\mathbb{R}$ , we have

$$d(a, c) + d(a', c') \leq d(a, c') \leq 4r.$$

Therefore either  $d(a, c) \leq 2r$  or  $d(a', c') \leq 2r$ , so either  $[a, c] \in \mathcal{V}$  or  $[a', c'] \in \mathcal{V}$ . But then either  $\gamma_\sigma$  or  $\gamma_{\sigma'}$  is nullhomotopic in  $\mathcal{V}$ , contradicting that  $G$  is a basis for  $\pi_1(\mathcal{V}, p)$ .  $\square$



**Figure 12.** Illustration of the argument by contradiction that  $|G| \leq 1$  in the proof of Theorem 6.13 (i). Critical simplices are black and matched simplices are gray. By the triangle inequality,  $d(a, c') \leq 4r$ , so since  $\{a < b < c \leq a' < b' < c'\}$  is isometric to a subset of  $\mathbb{R}$ , either  $[a, c] \in \mathcal{V}$  or  $[a', c'] \in \mathcal{V}$ .

*Proof of Theorem 6.13 (ii).* As in the statement of the theorem, let  $C_r$  denote the component of  $\mathcal{V}(P)_r$  containing  $\{p\}$ . We need to show that for  $r \leq r' \in [0, \infty)$ , if  $C_r \simeq S^1 \simeq C_{r'}$ , then the inclusion  $C_r \hookrightarrow C_{r'}$  is a homotopy equivalence. Let  $\gamma_\sigma : S^1 \rightarrow C_r$  and  $\gamma_{\sigma'} : S^1 \rightarrow C_{r'}$  be the generators for  $\pi_1(C_r, p)$  and  $\pi_1(C_{r'}, p)$  specified in the proof of Theorem 6.13 (i) above.

Given the way  $\gamma_\sigma$  and  $\gamma_{\sigma'}$  are defined, exactly one of the following must be true:

1.  $c \leq a'$ ,

2.  $c' \leq a$ ,
3.  $a \leq a' < c' \leq c$ .

We show that we cannot have  $c \leq a'$  using essentially the same triangle inequality argument we used in the proof of Theorem 6.13 (i): Suppose otherwise. Then  $d(a, p) < 2r$  and  $d(c', p) < 2r'$ . By the triangle inequality,  $d(a, c') \leq 2(r + r')$ , so we have

$$d(a, c) + d(a', c') \leq d(a, c') \leq 2(r + r').$$

Therefore either  $d(a, c) \leq 2r$  or  $d(a', c') \leq 2r'$ , leading to a contradiction as above.

The same argument shows that we cannot have  $c' \leq a$ . Therefore, we must have  $a \leq a' < c' \leq c$ .

We will show that if  $a \neq a'$ , then  $[a, a'] \in C_{r'}$ : We have  $d(a, p) \leq 2r$  and  $d(c', p) \leq 2r'$ , so by the triangle inequality,  $d(a, c') \leq 2(r + r')$ . Therefore either  $d(a, a') \leq 2r'$  or  $d(a', c') \leq 2r$ . But since  $\gamma_{\sigma'}$  is not nullhomotopic by assumption, we must have  $d(a', c') > 2r' \geq 2r$ , so  $d(a, a') \leq 2r'$ . Thus  $[a, a'] \in C_{r'}$ , as desired. It follows that  $[p, a, a'] \in C_{r'}$ .

The symmetric argument shows that if  $c' \neq c$ , then  $[p, c', c] \in C_{r'}$ . Letting

$$j : C_r \hookrightarrow C_{r'}$$

denote the inclusion, we thus have that  $j \circ \gamma_\sigma \sim \gamma_{\sigma'}$ . Since  $\gamma_\sigma$  and  $\gamma_{\sigma'}$  are both homotopy equivalences,  $j$  must be a homotopy equivalence as well.  $\square$

*Proof of Theorem 6.13 (iii).* Given the form of the set  $G$  of generators for  $\pi_1(\mathcal{V}(P)_r, p)$  given in the proof of Theorem 6.13 (i), it is clear that if

$$r \notin [d(p, P \setminus \{p\}), \text{diameter}(P \setminus \{p\})/2],$$

then  $\pi_1(\mathcal{V}(P)_r, p)$  is trivial. By (i) then, each component of  $\mathcal{V}(P)_r$  is contractible. Hence, the unique interval of  $\mathcal{B}_1(P)$ , if it exists, is contained in

$$[d(p, P \setminus \{p\}), \text{diameter}(P \setminus \{p\})/2].$$

To finish the proof of (iii), we need to show that the unique bar of  $\mathcal{B}_1(P)$  is of length at most  $d(p, P \setminus \{p\})$ . This follows from the stability of persistent homology. To see this, note that since  $P \setminus \{p\}$  is isometric to a subset of  $\mathbb{R}$ , it is tree-like, so Proposition 6.3 gives that  $\mathcal{B}_1(P \setminus \{p\}) = \emptyset$ . Therefore, by Theorem 5.5,

$$2d_B(\mathcal{B}_1(P), \emptyset) = 2d_B(\mathcal{B}_1(P), \mathcal{B}_1(P \setminus \{p\})) \leq d_H(P, P \setminus \{p\}) = d(p, P \setminus \{p\}),$$

where the last equality follows from the definition of  $d_H$ . The bottleneck distance of any barcode  $\mathcal{B}$  to the empty barcode is half the length of the longest interval of  $\mathcal{B}$ , so the result follows.  $\square$

**6.4. Inference about Recombination from Barcodes.** As an immediate corollary of the results of Sections 6.2 and 6.3, we now obtain our main result relating barcodes to recombination in the galled tree setting.

Recall from Section 3 that  $\mathcal{T}(\mathcal{E})$  denotes the topological novelty profile of a history  $\mathcal{E}$ , and that the temporal novelty of a recombinant  $r$  (with respect to some choice of time function) is denoted as  $\mathcal{N}(r)$ . Recall also that when  $\mathcal{E}$  is indexed by a galled tree,  $\mathcal{T}(\mathcal{E})$  is equal to the temporal novelty profile of  $\mathcal{E}$ , with respect to any time function.

For  $G$  a phylogenetic graph, let  $\mathcal{R}^G$  denote the set of recombinants of  $G$ . For  $\mathcal{B}$  a barcode, let  $\text{lengths}(\mathcal{B})$  denote the list of lengths of intervals of  $\mathcal{B}$ , sorted in descending order.



**Theorem 6.18.** *Let  $\mathcal{E}$  be a history indexed by a galled tree  $G$ .*

*(i) Theorem 6.11 (ii) and Theorem 6.13 (ii) yield a canonical injection*

$$\phi : \mathcal{B}_1(\mathcal{E}) \hookrightarrow \mathcal{R}^G,$$

*such that  $\text{length}(I) \leq \mathcal{N}(\phi(I))$  for all  $I \in \mathcal{B}_1(\mathcal{E})$ . In particular,*

$$\text{lengths}(\mathcal{B}_1(\mathcal{E})) \leq \mathcal{T}(\mathcal{E}).$$

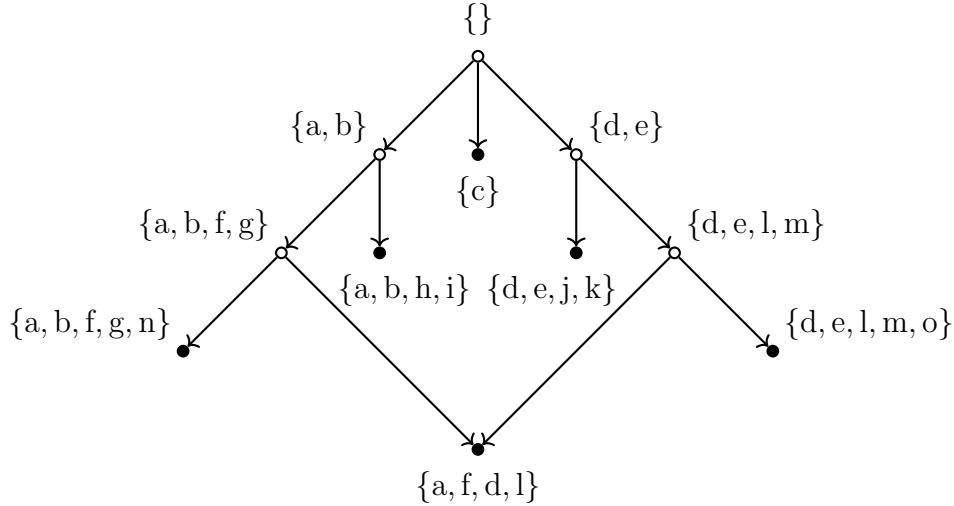
*(ii)  $\mathcal{B}_i(\mathcal{E}) = \emptyset$  for  $i \geq 2$ .*

*Proof.* For  $G$  a galled tree, each  $r \in \mathcal{R}^G$  corresponds to an entry of  $\mathcal{T}(\mathcal{E})$ ; in fact, this entry is easily seen to be  $d(\mathcal{E}^{L \setminus r}, \mathcal{E}_r)$ , where  $L$  denotes the directed loop corresponding to  $R$ , and  $\mathcal{E}^{L \setminus r}$  denotes the restriction of  $\mathcal{E}$  to vertices of  $L$  other than  $r$ . (i) now follows from Theorem 6.13 (iii).

(ii) is immediate from Theorem 6.11 (ii), Proposition 6.12, and Theorem 6.13 (i).  $\square$

**Example 6.19.** Given the analogy between Theorem 6.4 (for trees) and Theorem 6.18 (for galled trees), and the fact that Theorem 6.4 holds for arbitrary subsamples of a history, it is natural to ask whether Theorem 6.18 also holds for arbitrary subsamples. The example shown in Fig. 13 demonstrates that Theorem 6.18 (i) does not hold for arbitrary subsamples; the example, discovered by computer, is a subset  $S$  of a history indexed by a galled tree with a single recombinant, for which  $\mathcal{B}_1(S) = \{[5, 6], [5, 6]\}$ . We conjecture that Theorem 6.18 (ii) also does not hold for arbitrary subsets.

Nevertheless, it may be the case that for reasonable random models of histories indexed by galled trees, violations of Theorem 6.18 are relatively rare. We provide some preliminary numerical evidence for this in Appendix B, focusing on how often the number of intervals in  $\mathcal{B}_1(S)$  of a sample  $S$  exceeds the number of recombinants in the underlying history.



**Figure 13.** A subset  $S$  of a history indexed by a galled tree with one recombinant, for which  $|\mathcal{B}_1(S)| = 2$ . Nodes corresponding to elements of  $S$  are shown in black; the remaining nodes are shown in white.

## 7. RELAXING THE COMPLETE SAMPLING AND GALLED TREE ASSUMPTIONS

Theorem 6.18, the main result of the previous section, holds under the assumption that our evolutionary history is indexed by a galled tree, and that all organisms in the history have been sampled. In this section, we apply the stability of persistent homology to extend the theorem to the case of an arbitrary (noisy) subsample of a history indexed by an arbitrary phylogenetic graph.

**7.1. Relaxing the Complete Sampling Assumption.** First, we extend Theorem 6.18 to the case of a noisy subsample. Given a list of non-negative numbers  $L$ , let  $\text{Trim}(L, \delta)$  be the list obtained by removing each of the numbers less than or equal to  $\delta$  and subtracting  $\delta$  from each of the remaining numbers.

**Corollary 7.1.** *Let  $\mathcal{E}$  be a history indexed by a galled tree and let  $S$  be a finite metric space with  $d_{GH}(\mathcal{E}, S) = \delta$ . Then*

- (i)  $\text{Trim}(\text{lengths}(\mathcal{B}_1(S)), 2\delta) \leq \mathcal{T}(\mathcal{E})$ .
- (ii) For  $i \geq 2$ , each interval of  $\mathcal{B}_i(S)$  has length at most  $2\delta$ .

*Proof.* This follows immediately from Theorems 5.4 and 6.18. □

**7.2. Relaxing the Galled Tree Assumption.** As an application of Corollary 7.1 (i), we next also relax the assumption that  $\mathcal{E}$  is indexed by a galled tree, yielding a further extension of Theorem 6.18 which applies to any phylogenetic graph.

For  $G$  any phylogenetic graph and  $\mathcal{E}$  a history indexed by  $G$ , let

$$\bigcup \mathcal{E} := \bigcup_{v \in V} \mathcal{E}_v.$$

Thus,  $\bigcup \mathcal{E}$  is set of all mutations appearing in the history  $\mathcal{E}$ .

For  $M \subset \bigcup \mathcal{E}$  any subset and  $v \in V$ , let  $\mathcal{E}_v^M = \mathcal{E}_v \setminus M$ . Let  $G^M$  denote a subgraph of  $G$  obtained by removing edges as follows: Suppose  $w$  is a recombinant of  $G$  with parents  $u, v$ . If  $\mathcal{E}_w^M = \mathcal{E}_u^M \neq \mathcal{E}_v^M$ , we remove the edge  $(v, w)$  from  $G$ . If  $\mathcal{E}_w^M = \mathcal{E}_u^M = \mathcal{E}_v^M$  we remove exactly one of the edges  $(u, w)$  and  $(v, w)$ , choosing arbitrarily. It is easy to check that the sets  $\mathcal{E}_v^M$  then give a well-defined evolutionary history  $\mathcal{E}^M$  indexed by  $G^M$ .

**Definition 7.2.** We let

$$\text{Gall}(\mathcal{E}) := \min \left\{ |M| \mid M \subset \bigcup \mathcal{E} \text{ such that } G^M \text{ can be chosen to be a galled tree} \right\}.$$

Informally,  $\text{Gall}(\mathcal{E})$  is the number of mutations in  $\mathcal{E}$  which must be ignored to obtain a history indexed by a galled tree by pruning edges in  $G$ .

The following is our most general result relating barcodes and recombination:

**Corollary 7.3.** *Let  $\mathcal{E}$  be a history indexed by an arbitrary phylogenetic graph  $G$ , and let  $S$  be a finite metric space with  $d_{GH}(\mathcal{E}, S) = \delta$ . Then*

- (i)  $\text{Trim}(\text{lengths}(\mathcal{B}_1(S)), 3 \text{Gall}(\mathcal{E}) + 2\delta) \leq \mathcal{T}(\mathcal{E})$ .
- (ii) For  $i \geq 2$ , each interval of  $\mathcal{B}_i(S)$  has length at most  $2(\text{Gall}(\mathcal{E}) + \delta)$ .

*Proof.* Choose  $M \subset \bigcup \mathcal{E}$  and a galled tree  $G^M$  as above, such that  $|M| = \text{Gall}(\mathcal{E})$ . Note that  $d_{GH}(\mathcal{E}^M, \mathcal{E}) \leq |M| = \text{Gall}(\mathcal{E})$ , so by the triangle inequality,  $d_{GH}(\mathcal{E}^M, S) \leq \text{Gall}(\mathcal{E}) + \delta$ . (ii) then follows from Corollary 7.1 (ii). By Corollary 7.1 (i),

$$\text{Trim}(\text{lengths}(\mathcal{B}_1(S)), 2(\text{Gall}(\mathcal{E}) + \delta)) \leq \mathcal{T}(\mathcal{E}^M). \quad (1)$$

Letting  $\bar{\mathcal{T}}(\mathcal{E}^M)$  be the vector of length  $|\mathcal{T}(\mathcal{E})|$  obtained by adding some 0's to the end of  $\mathcal{T}(\mathcal{E}^M)$ , we have by Proposition 5.6 that  $d_\infty(\bar{\mathcal{T}}(\mathcal{E}^M), \mathcal{T}(\mathcal{E})) \leq |M|$ . Together with (1), this implies that

$$\text{Trim}(\text{lengths}(\mathcal{B}_1(S)), 3 \text{Gall}(\mathcal{E}) + 2\delta) \leq \mathcal{T}(\mathcal{E}),$$

which gives (i).  $\square$

**Remark 7.4.** Clearly, for Corollary 7.3 to yield a strong bound,  $\text{Gall}(\mathcal{E})$  must be small. One might expect  $\text{Gall}(\mathcal{E})$  to be small but non-zero when recombination events typically affect short genome tracts (e.g., when they are gene conversion events).

## 8. RANDOM HISTORIES INDEXED BY GALLED TREES

The results we have presented so far have been deterministic. In Section 8.1 below, we observe that in a wide class of probabilistic models of genetic sequence evolution on galled trees, the intervals of the first persistence barcode are independent random variables. Thus, to understand the statistical properties of these barcodes, it suffices to understand the special case that the galled tree is a directed loop. In Section 8.2, we study this special case numerically, for one choice of probabilistic model.

**8.1. Independence of Intervals in the First Barcode.** In this section, we assume the reader is familiar with basic elements of the modern, measure-theoretic formulation of probability theory [21] and with the definition of conditional independence given a random variable [36].

*Notation.* Suppose  $X$ ,  $Y$ , and  $Z$  are random variables on the same probability space. If  $X$  is independent of  $Y$ , we write  $X \perp\!\!\!\perp Y$ . If  $X$  is independent of  $Y$  given  $Z$ , we write  $X \perp\!\!\!\perp Y \mid Z$ .

For  $P$  a poset and  $p \in P$ , let

$$\text{nd}(p) := \{q \in P \mid p \not\leq q\}.$$

(Here,  $\text{nd}$  stands for *non-descendants*.) In what follows,  $P$  will often be the vertex set of a directed acyclic graph, with the partial order induced by the graph.

If  $\mathcal{E}$  is an evolutionary history indexed by  $G$ ,  $V$  is the vertex set of  $G$ , and  $S \subset V$ , we write  $\mathcal{E}_S := \{\mathcal{E}_v \mid v \in S\}$ . Similarly, for  $v, w \in V$ , let  $\mathcal{E}_{v \setminus w} := \mathcal{E}_v \setminus \mathcal{E}_w$ , and  $\mathcal{E}_{v \cap w} := \mathcal{E}_v \cap \mathcal{E}_w$ . We will also use these notation conventions in combination with one another, so that e.g.,  $\mathcal{E}_{r \setminus (p \cap q)}$  is understood to denote  $\mathcal{E}_r \setminus (\mathcal{E}_p \cap \mathcal{E}_q)$ .

**Definition 8.1** (Random History). For  $G$  a fixed phylogenetic graph with vertices  $V$ , a *random (evolutionary) history*  $\mathcal{E}$  indexed by  $G$  consists of the following data:

- A probability space  $\Omega$ .
- A countable set  $X_v$  for each  $v \in V$ , such that each element of  $X_v$  is itself a set. We equip  $X_v$  with the discrete  $\sigma$ -algebra.

- For each  $v \in V$ , a random variable  $\mathcal{E}_v : \Omega \rightarrow X_v$  such that for each  $\omega \in \Omega$ ,  $\{\mathcal{E}_v(\omega)\}_{v \in V}$  is an evolutionary history.

**Definition 8.2** (Locally Markov History). Suppose that  $\mathcal{E}$  is a random history indexed by a phylogenetic graph  $G$  with vertices  $V$ .  $\mathcal{E}$  is said to be *locally Markov* if for each  $v \in V$ ,

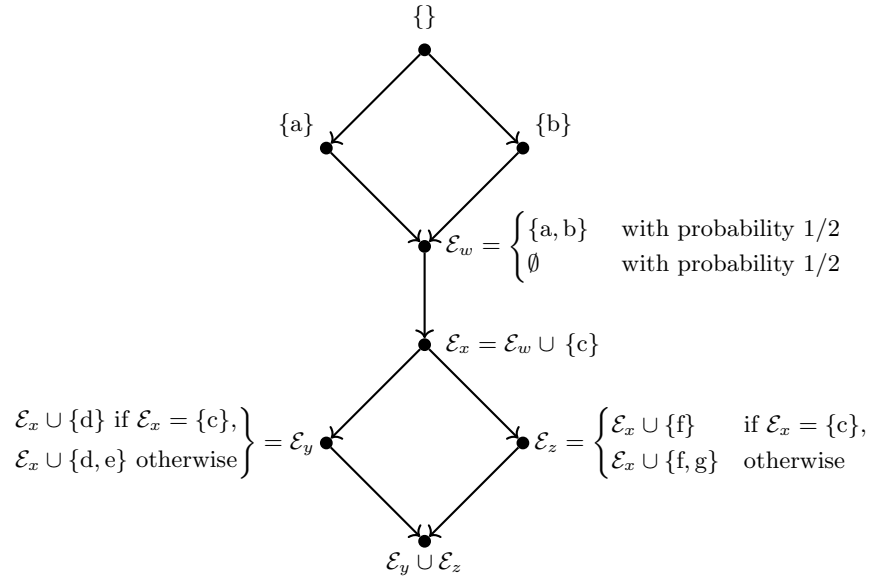
$$\mathcal{E}_v \perp\!\!\!\perp \mathcal{E}_{\text{nd}(v)} \mid \{\mathcal{E}_p \mid p \text{ a parent of } v\}. \quad (2)$$

A locally Markov history is a special case of a *Bayesian network*, a widely used probabilistic model [38].

The assumption that a random history is locally Markov is quite natural; informally, this says that the genome of each organism depends only on the genomes of its parents. However, the next example shows that for  $\mathcal{E}$  a locally Markov history, it is not necessarily the case that the intervals in  $\mathcal{B}_1(\mathcal{E})$  are independent.

**Example 8.3.** In the locally Markov history  $\mathcal{E}$  of Fig. 14, the mutations from the top directed loop are passed down to the bottom directed loop, where they serve as “instructions” for how clonal mutations occur in the bottom loop. Thus, the intervals in  $\mathcal{B}_1(\mathcal{E})$  associated to the two recombinants are not independent.

For each vertex  $v \neq w$ ,  $\mathcal{E}_v$  is completely determined by its parents. The top recombinant corresponds to an interval  $[\frac{1}{2}, 1)$  in the first barcode with probability  $\frac{1}{2}$ , and to an empty interval with probability  $\frac{1}{2}$ . If the top recombinant corresponds to  $[\frac{1}{2}, 1)$ , then the bottom recombinant corresponds to  $[1, 2)$ ; otherwise, the bottom recombinant corresponds to  $[\frac{1}{2}, 1)$ .



**Figure 14.** A locally Markov history  $\mathcal{E}$  for which the intervals in the 1<sup>st</sup> persistence barcode corresponding to the two recombinants are not independent.

Motivated by the above, we introduce the following subclass of locally Markov histories:

**Definition 8.4** (Phylogenetically Markov History). Suppose that  $\mathcal{E}$  is a random history indexed by a fixed phylogenetic graph  $G$  with vertices  $V$ . We say  $\mathcal{E}$  is *phylogenetically Markov* if

1.  $\mathcal{E}$  is locally Markov.
2.  $\mathcal{E}_{v \setminus p} \perp\!\!\!\perp \mathcal{E}_p$  for all clones  $v \in V$  with parent  $p$ .
3. For all recombinants  $r$  with parents  $p$  and  $q$ ,

$$\mathcal{E}_{r \setminus (p \cap q)} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(r)} \mid \mathcal{E}_{\{p \setminus (p \cap q), q \setminus (p \cap q)\}}.$$

To parse condition 3, recall that a recombinant  $r$ 's genome necessarily inherits what is common to both parents  $p$  and  $q$ ; condition 3 states that the rest of  $r$ 's genome is independent of the genomes of all non-descendants of  $r$ , given the rest of the genomes of each parent.

**Remark 8.5.** Definition 8.4 is slightly redundant, in the sense that condition 3 implies the local Markov property for each recombinant; this follows from Lemma 8.8 below. One might hope that one could obtain an equivalent definition by replacing condition 3 in Definition 8.4 with the simpler condition that  $\mathcal{E}_{r \setminus (p \cap q)} \perp\!\!\!\perp \mathcal{E}_{p \cap q}$ , but in fact this is strictly weaker. It can be shown that our independence result (Theorem 8.7 below) does not hold for this weaker condition.

**Example 8.6.** Assume that  $G$  is endowed with a time function  $t : V \rightarrow \mathbb{R}$ , as defined in Section 3. We specify (up to choice of labels for mutations) a phylogenetically Markov history  $\mathcal{E}$ , the *Poisson history* indexed by  $G$ :

- $\mathcal{E}_r = \emptyset$ , for  $r$  the root of  $G$ .
- If  $w$  is a clone with ancestor  $v$ ,  $|\mathcal{E}_w \setminus \mathcal{E}_v|$  is Poisson distributed with parameter  $t(w) - t(v)$ .
- If  $w$  is a recombinant with parents  $u$  and  $v$ , then for each  $m \in \mathcal{E}_u \setminus \mathcal{E}_v$ ,  $P(m \in \mathcal{E}_w) = p_w$ , and for each  $m \in \mathcal{E}_v \setminus \mathcal{E}_u$ ,  $P(m \in \mathcal{E}_w) = 1 - p_w$ . Here, we may either take  $p_w = 1/2$  for all  $w$ , or take the  $p_w$  to be i.i.d. random variables with the uniform distribution on  $[0, 1]$ .

Let

$$\mathcal{I} := \{[a, b) \mid a < b \in \{0, 1, 2, \dots\}\} \cup \{\emptyset\}.$$

Thus,  $\mathcal{I}$  is a collection of intervals with integer endpoints, together with the empty interval. For  $\mathcal{E}$  a history indexed by a galled tree  $G$  and  $r$  a recombinant in  $G$ , let  $\mathcal{I}^\mathcal{E}(r) \in \mathcal{I}$  denote the unique interval in  $\mathcal{B}_1(\mathcal{E})$  corresponding to  $r$ , if such an interval exists (see Theorem 6.18), and let  $\mathcal{I}^\mathcal{E}(r) = \emptyset$  otherwise. As in Section 6.4, we let  $\mathcal{R}^G$  denote the set of recombinants of  $G$ .

Here is the main result of this section:

**Theorem 8.7.** *For  $\mathcal{E}$  a phylogenetically Markov history indexed by a galled tree  $G$ , the random variables  $\{\mathcal{I}^\mathcal{E}(r)\}_{r \in \mathcal{R}^G}$  are independent.*

The proof of the theorem will use several standard facts about conditional independence, which we record in the following lemma.

**Lemma 8.8.** *Assume  $h$  is a measurable function whose domain is the codomain of the random variable  $X$ .*

- (i) If  $X \perp\!\!\!\perp Y \mid Z$ , then  $Y \perp\!\!\!\perp X \mid Z$ .
- (ii) If  $X \perp\!\!\!\perp Y \mid Z$ , then  $h(X) \perp\!\!\!\perp Y \mid Z$ .
- (iii) If  $X \perp\!\!\!\perp Y \mid Z$ , then  $X \perp\!\!\!\perp Y \mid (Z, h(X))$ .
- (iv) If  $X \perp\!\!\!\perp Y \mid Z$  and  $W \perp\!\!\!\perp Y \mid (X, Z)$ , then  $(W, X) \perp\!\!\!\perp Y \mid Z$ .
- (v) If  $X \perp\!\!\!\perp Y \mid Z$ , then  $(Z, X) \perp\!\!\!\perp Y \mid Z$ .

Note that by taking  $Z$  to be the identity random variable, we also obtain unconditional versions of (i)-(iv) above.

*Proof.* Properties (i)-(iv) appear in many places; see e.g. [38]. We prove (v). Taking  $W = Z$  in (iv), it suffices to show that  $Z \perp\!\!\!\perp Y \mid (X, Z)$ . By (ii), for this it is enough to show that  $(X, Z) \perp\!\!\!\perp Y \mid (X, Z)$ . But it is easy to check that in general,  $A \perp\!\!\!\perp B \mid A$ .  $\square$

**Lemma 8.9.** *If  $\{X_a\}_{a \in P}$  is a collection of random variables indexed by a finite poset  $P$  and  $X_a \perp\!\!\!\perp X_{\text{nd}(a)}$  for each  $a \in P$ , then the  $\{X_a\}_{a \in P}$  are independent.*

*Proof.* Choose a total order compatible with the partial order on  $P$ , and relabel the random variables with respect to this order as  $X_1, \dots, X_{|P|}$ . We show by induction that  $X_1, \dots, X_m$  are independent for each  $m \in \{1, \dots, |P|\}$ . The base case is trivial. Now suppose  $X_1, \dots, X_{m-1}$  are independent. The elements of  $P$  corresponding to the indices  $1, \dots, m-1$  are in  $\text{nd}(m)$ , so  $X_m$  is independent of  $X_1, \dots, X_{m-1}$ . By this and the induction hypothesis,  $X_1, \dots, X_m$  are independent.  $\square$

*Proof of Theorem 8.7.* Order the vertices of  $G$  arbitrarily. For  $r \in \mathcal{R}^G$ , let  $D^r$  denote the distance matrix obtained by restricting  $\mathcal{E}$  to the directed loop  $L^r$  of  $G$  with sink  $r$ . The images of independent random variables under measurable functions remain independent, so in view of the results of Section 6, it suffices to show that the  $\{D^r\}_{r \in \mathcal{R}^G}$  are independent. Now for each recombinant  $r$ , let  $V^r$  denote the vertices of  $L^r$ , and let  $q_r$  denote the unique source of  $L^r$ . Since  $\mathcal{E}_{q_r} \subset \mathcal{E}_v$  for all  $v \in V^r$ , clearly  $D^r$  is determined by  $A^r := \{\mathcal{E}_{v \setminus q_r}\}_{v \in V^r \setminus \{q_r\}}$ . Therefore, it in fact suffices to show that the random sets  $\{A^r\}_{r \in \mathcal{R}^G}$  are independent.

We define a partial order on  $\mathcal{R}^G$  by writing  $r \leq r'$  if for some  $v \in V^r \setminus \{q_r\}$ , there is a directed path from  $v$  to  $q_{r'}$  in  $G$ ; it is easy to check that this is in fact a partial order. This partial order induces a partial order on  $\{A^r\}_{r \in \mathcal{R}^G}$ . We establish the independence of the  $\{A^r\}_{r \in \mathcal{R}^G}$  by applying Lemma 8.9, using this partial order. Let

$$\text{nd}(L^r) := \bigcap_{v \in V^r \setminus \{q_r\}} \text{nd}(v).$$

If  $r \not\leq r'$  then  $V^{r'} \subset \text{nd}(L^r)$ , so  $\{A^{r'} \mid A^r \not\leq A^{r'}\}$  is completely determined by  $\mathcal{E}_{\text{nd}(L^r)}$ . Thus, it suffices to show that  $A^r \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$  for each  $r \in \mathcal{R}^G$ .

Let us fix  $r \in \mathcal{R}^G$  and write  $q = q^r$ . Choose a total order on  $V^r \setminus \{r, q\}$  compatible with the partial order on  $V$ , and write the elements in increasing order as  $\{c_1, \dots, c_m\}$ . For  $j \in \{1, \dots, m\}$ , let  $B_j := \{\mathcal{E}_{c_i \setminus q}\}_{1 \leq i \leq j} \subset A^r$ . We show by induction on  $j$  that  $B_j \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$  for each  $j$ .

First, consider the base case  $j = 1$ . In the remainder of the proof, the five statements of Lemma 8.8 will be denoted simply as (i)-(v). By the definition of a phylogenetically Markov history, we have  $\mathcal{E}_{c_1} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(c_1)} \mid \mathcal{E}_q$ , so by (ii), we have  $\mathcal{E}_{c_1} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)} \mid \mathcal{E}_q$ . By (v) then,  $\mathcal{E}_{\{c_1, q\}} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)} \mid \mathcal{E}_q$ , so by (ii),  $\mathcal{E}_{c_1 \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)} \mid \mathcal{E}_q$ . The definition of a phylogenetically

Markov history also gives that  $\mathcal{E}_{c_1 \setminus q} \perp\!\!\!\perp \mathcal{E}_q$ . Applying (iv) and (ii), we find that  $\mathcal{E}_{c_1 \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ . This shows that  $B_1 \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ .

The induction step is similar to the above. Let  $p$  denote the parent of  $c_j$ .  $\mathcal{E}_{c_j} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(c_j)} \mid \mathcal{E}_p$ , so  $\mathcal{E}_{c_j \setminus p} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(c_j)} \mid \mathcal{E}_p$ . Moreover,  $\mathcal{E}_{c_j \setminus p} \perp\!\!\!\perp \mathcal{E}_p$ , so  $\mathcal{E}_{c_j \setminus p} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(c_j)}$ . Then by (ii) and (iii),  $\mathcal{E}_{c_j \setminus p} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)} \mid B_{j-1}$ . By (iv) and the induction hypothesis, we thus have that  $B_j \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ , as desired.

Finally, we show that  $A^r \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ . Let  $p_1$  and  $p_2$  denote the parents of  $r$ . By the third condition in the definition of a phylogenetically Markov history, we have  $\mathcal{E}_{r \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(r)} \mid \mathcal{E}_{\{p_1 \setminus q, p_2 \setminus q\}}$ . By (iii),  $\mathcal{E}_{r \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(r)} \mid (\mathcal{E}_{\{p_1 \setminus q, p_2 \setminus q\}}, B_m)$ .  $\sigma(\mathcal{E}_{\{p_1 \setminus q, p_2 \setminus q\}}, B_m) = \sigma(B_m)$  since  $\mathcal{E}_{\{p_1 \setminus q, p_2 \setminus q\}} = h(B_m)$  for some measurable function  $h$ , so  $\mathcal{E}_{r \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(r)} \mid B_m$  by the definition of conditional independence. By (ii),  $\mathcal{E}_{r \setminus q} \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)} \mid B_m$ . We have also shown that  $B_m \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ , so by (iv),  $A^r \perp\!\!\!\perp \mathcal{E}_{\text{nd}(L^r)}$ .  $\square$

## 8.2. The Barcode of a Random History on a Directed Loop: Numerical Results.

Theorem 8.7 tells us that for a phylogenetically Markov history indexed by a galled tree, to understand the distribution of the 1<sup>st</sup> barcode, it suffices to understand this for each directed loop in the galled tree. Working with a simple random model of a history  $\mathcal{E}$  indexed by a directed loop, we now study the distribution of  $\mathcal{B}_1(\mathcal{E})$  numerically. Recall that by Theorem 6.13 (ii),  $\mathcal{B}_1(\mathcal{E})$  has at most one interval. We consider here the probability that  $\mathcal{B}_1(\mathcal{E})$  is nontrivial, as well as the average length of the interval.

In our simulations, we find that in the limit of high novelty, persistent homology captures between 14% and 35% of recombination events, the exact value depending on mutational parameters. In typical simulations where a recombination event is detected, the bar length is well below the theoretical maximum provided by Theorem 6.18 (i), scaling roughly as the square root of novelty.

*Details of the Computations.* We now specify the random model of a history indexed by a directed loop that we use in our simulations. The model depends on parameters  $m$  and  $k$ . Each random history generated by this model consists of: a left parent with no mutations; a right parent with mutations  $\{1, \dots, m\}$ ; a recombinant with some subset of these mutations; and  $k$  “intermediate sequences,” each randomly sampled (with replacement) from the set

$$\{\{1\}, \{1, 2\}, \{1, 2, 3\}, \dots, \{1, \dots, m-1\}\}.$$

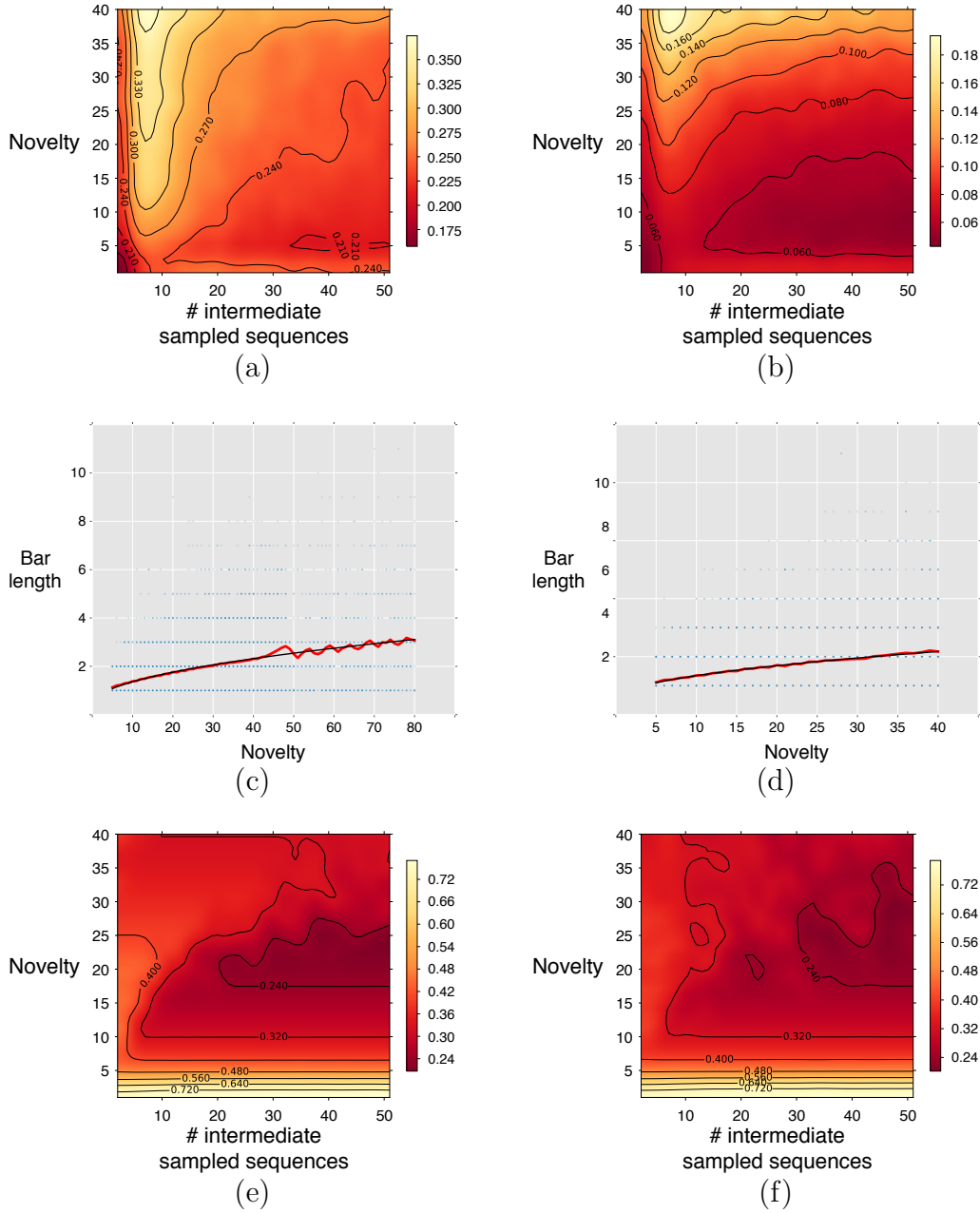
In our simulations, we consider values of  $m$  between 2 and 200, and values of  $k$  between 1 and 50. In addition, we consider a “maximal sampling” scenario, in which all possible intermediate sequences were included in the sample; for the purpose of visualization, this scenario is assigned parameter value  $k = 51$ . To construct the recombinant, we select each of the  $m$  mutations with probability  $\alpha$ . In one set of our simulations, we set  $\alpha = 0.5$  (simulating a recombination breakpoint at the midpoint of the genome); in a second set of simulations, we choose  $\alpha$  randomly from the uniform distribution on  $[0, 1]$ . For each sampled history, we compute both the novelty of the recombinant and the persistent homology of the sample.

*Results.* The results of our simulations are given in figure Fig. 15. To obtain each subfigure, we aggregated the data for the various values of the parameter  $m$ . We see that the rate of detection of a recombinant increases with novelty, up to about 37% (midpoint recombination

breakpoint, Fig. 15 (a)) or 20% (uniform recombination breakpoint, Fig. 15 (b)). For high novelty, increasing  $k$  improves detection only up to about  $k = 7$ , after which detection falls to about 28% (midpoint recombination breakpoint) or 16% (uniform recombination breakpoint) for high  $k$ .

Bar length typically falls well below the upper bound given by the novelty of the recombinant. In particular, for simulations where the recombination event was detected, bar length scales roughly as the square root of novelty (Fig. 15 (c,d)). For cases with high novelty, median bar length ranges from about 25% of the square root of novelty (if many intermediate sequences are sampled, upper right corner of Fig. 15 (e,f)) to 35% of the square root of novelty (if few intermediate sequences are sampled, upper left corner of Fig. 15 (e,f)).





**Figure 15.** Sensitivity of persistent homology in simulations of single recombination events. Panels (a,b): Fraction of simulations in which the recombination event was detected. Recombination breakpoint is either at the midpoint of (panel a) or uniformly distributed along (panel b) the genome. Panels (c,d): Each translucent point marks the result of a single simulation in which the recombination event was detected, the red line tracks the average bar length among all simulations with equal novelty, and the black line shows the least-squares fit of parameters  $a$  and  $b$  among functions  $y = a\sqrt{x} + b$ . Panel (c): Recombination breakpoint at midpoint of the genome. The fit of all 684,026 cases is  $y = 0.30 \times \sqrt{x} + 0.40$ . Panel (d): Recombination breakpoint uniformly distributed along the genome. The fit of all 134,830 cases is  $y = 0.26 \times \sqrt{x} + 0.51$ . Panels (e,f): Median of ratio of bar length to square root of novelty, conditional on the recombination event being detected (i.e., bar length  $\geq 1$ ). Recombination breakpoint is either at the midpoint of (panel e) or uniformly distributed along (panel f) the genome.

## 9. DISCUSSION

In this paper, we have introduced *novelty profiles*, simple statistics of an evolutionary history which not only count the number of recombination events in the history, but also quantify the contribution recombination makes to genetic diversity in the population. We have studied the problem of inferring information about a novelty profile from the persistent homology of sampled data. Our results provide mathematical foundations for several earlier works which have used persistent homology to study recombination.

*Potential Applications of the Novelty Profile.* It is generally accepted that genetic recombination – including horizontal gene transfer in bacteria and sexual reproduction in eukarya – evolved both to hasten adaptive evolution [18, 39] and to prevent “mutational meltdown” caused by the accumulation of deleterious mutations [27]. Understanding the precise mechanisms by which recombination contributes to evolution is a long-standing problem in evolutionary biology. An important subproblem is to understand how recombination shapes the genomic distributions on which natural selection acts. As a population’s genomic distribution is affected not only by rates of mutation and recombination but also by the genetic novelty of recombinants, the novelty profile may be a useful tool for addressing this subproblem.

We describe one class of potential applications in this direction. It is well known that recombination plays an important role both in the spread of infection and in the development of drug resistance. For example: reassortment can cause outbreaks of influenza (e.g., the Swine flu pandemic of 2009) [34, 45, 49]; recombination can effect resistance to anti-viral medication in HIV [42]; and horizontal gene transfer plays a major role in the spread of antibiotic resistance in bacteria [20]. The novelty profile could be useful for developing a fuller quantitative understanding of the role of recombination in such epidemiological events. We hypothesize that, compared to a count of recombination events alone, the novelty profile better predicts both future outbreaks of infection and the proliferation of drug resistance. It could be interesting to test this hypothesis in simulation. To test the hypothesis on real biological data, one needs well-behaved estimators of the novelty profile; our results represent progress in this direction. If the novelty profile is indeed predictive of outbreaks or of the proliferation of drug resistance, statistics derived from estimates of novelty profiles could potentially inform public health responses to infectious disease.

*On the Assumptions Underlying our Main Results.* Our main results relating barcodes to novelty profiles depend on strong assumptions about the evolving population and genomic sampling. In their simplest form, our results assume that the evolutionary history  $\mathcal{E}$  is indexed by a galled tree and that all genomes in the history are included in our sample. Using the stability of persistent homology, we have extended these results to hold for an arbitrary sample  $S$  of an arbitrary evolutionary history  $\mathcal{E}$ . The strength of the bounds provided by these extended results, relative to the ideal case of a galled tree with every genome sampled, is controlled by  $d_{GH}(S, \mathcal{E})$ , the Gromov-Hausdorff distance between  $S$  and  $\mathcal{E}$ , and  $\text{Gall}(\mathcal{E})$ , the number of mutations in  $\mathcal{E}$  which must be ignored to obtain a history indexed by a galled tree. In cases where  $d_{GH}(S, \mathcal{E})$  and  $\text{Gall}(\mathcal{E})$  can be assumed to be small relative to the lengths of intervals in the barcodes  $\mathcal{B}_i(S)$ , our results provide an informative lower bound on the novelty

profile, though the numerical results of Section 8.2 suggest that this bound is typically far from tight.

These results raise three key questions about applications of our work: First, under what circumstances can real-world genomic samples be expected to exhibit small enough values of  $d_{GH}(S, \mathcal{E})$  and  $\text{Gall}(\mathcal{E})$  for our results to yield a useful bound on the novelty profile? Second, can how can the theory developed in this paper be extended to yield a useful topological bound on the novelty profile in cases where cannot expect  $d_{GH}(S, \mathcal{E})$  or  $\text{Gall}(\mathcal{E})$  to be small? And third, can more sensitive bounds on the novelty profile be obtained? We discuss each of these questions in turn.

*The Small  $\text{Gall}(\mathcal{E})$  Condition.* While restrictive, the assumption that  $\text{Gall}(\mathcal{E})$  is small is biologically plausible in some settings. As shown in Appendix A in the context of the coalescent model, the phylogenetic graph  $G$  indexing a history  $\mathcal{E}$  will be a galled tree with high probability if and only if a relatively strong condition on the rareness of recombination is satisfied. It is important to note that this condition depends not only on the species studied, but also the sample size and how the genomic data is analyzed. For example, in the study of human recombination, a key methodological choice is the size of the genomic window used for the analysis; while a larger window offers more accurate estimates of recombination rate, a smaller window better localizes recombination breakpoints [9]. Since recombination is also rarer in a smaller window, a sample's ancestry is more likely to be represented by a galled tree when using a smaller window. In a sample of 125 humans, for instance, the ancestry in an average 275 bp window is predicted to be a galled tree with 90% probability (Appendix A). As recombination rates vary dramatically across segments of the human genome, the actual probability likewise varies.

Theorem 6.18 guarantees that histories indexed by galled trees have empty barcodes  $\mathcal{B}_i(\mathcal{E})$  for  $i \geq 2$ . As such, applications of our work may involve choosing a window size for which this condition holds. However, empty barcodes  $\mathcal{B}_i(\mathcal{E})$  for  $i \geq 2$  do not guarantee that  $\mathcal{E}$  is indexed by a galled tree; it remains to better understand how useful higher barcodes are in practice for determining whether  $\text{Gall}(\mathcal{E})$  is small.

*Extending our Results to More Complex Phylogenetic Graphs.* It may be possible to extend our results to histories indexed by iterated sums of phylogenetic graphs with at most  $k$  recombination events, at least for small  $k$ . (The galled tree setting is the case  $k = 1$ .) Such an extension would then yield informative lower bounds on the novelty profile for a larger class of evolutionary histories.

The next logical step would be to study the case  $k = 2$ . In analogy with our main results, several questions arise about a history  $\mathcal{E}$  indexed by a phylogenetic graph with two recombinants:

- Does  $\mathcal{B}_1(\mathcal{E})$  have at most two bars?
- For which degrees  $i$  is  $\mathcal{B}_i(\mathcal{E})$  necessarily trivial?
- What can be said about the lengths of the intervals of  $\mathcal{B}_i(\mathcal{E})$ ?

We have conjectured (Remark 6.15) that our results about almost linear metric spaces extend to almost tree-like metric spaces. Proving this may be a useful first step towards answering the above questions.

*The Small Gromov-Hausdorff Distance Condition.* For a typical genomic sample  $S$ ,  $d_{GH}(S, \mathcal{E})$  can be large. Indeed, regardless of whether individuals are sampled simultaneously or longitudinally, the most recent common ancestor of two individuals in  $S$  may be genetically distant from all individuals in  $S$ .

There are applications, however, where sampling is so dense, and so frequent, that we do expect common ancestors of sampled individuals to be genetically close to individuals in the sample. One such application is dense epidemiological sampling of HIV, for which entire countries have established long-term viral genomic surveillance [6, 48]. In such cases, standard phylogenetic methods, which assume that ancestors are absent from the sample, may produce misleading results [47]. As genetic sequencing continues to decline in cost, it is likely for dense longitudinal genomic samples to become more common. For such data sets, particularly for pathogens where evolutionary time scales are short compared to sampling duration, the assumption that  $d_{GH}(S, \mathcal{E})$  is small may be more reasonable.

Even for samples  $S$  for which  $d_{GH}(S, \mathcal{E})$  is large, our simulation results using the coalescent model in Appendix B suggest that violations of the exact bound on the number of recombinations given by Theorem 6.18 (i) are relatively rare, in part because of the limited sensitivity of persistence barcodes in detecting recombination (Section 8.2). This, together with the empirical results from extensive simulations described in previous literature [8, 9, 12, 23], give us hope that our main theoretical results may be extended to probabilistic ones that yield useful bounds even when  $d_{GH}(S, \mathcal{E})$  is large.

*Tightening our Lower Bounds on the Novelty Profile.* As shown numerically in Section 8.2, one limitation of barcodes as estimators of novelty profiles is their relatively low sensitivity to individual recombination events. A natural goal is to devise a more sensitive variant of the barcode estimator, with similar theoretical guarantees. Is it possible to develop a *consistent* barcode estimator for the novelty profile, for a reasonable class of probabilistic models?

As a step in this direction, it would be interesting to apply our independence result, Theorem 8.7, to obtain analytic results about the probability distribution on  $\mathcal{B}_1(\mathcal{E})$ , for  $\mathcal{E}$  a phylogenetically Markov random evolutionary history (e.g., Poisson) indexed by a galled tree. Ideally, such results would explain the relationships between novelty and the barcode observed empirically in Section 8.2.

*Estimating Novelty Profiles by Other Means.* While this paper has focused primarily on bounding the novelty profile using persistent homology, other approaches to estimating the novelty profile may be fruitful as well. For example, there is a large literature on direct estimation of evolutionary histories (which, as noted earlier, are usually called *ancestral recombination graphs*); see for example [28] and the references therein. Though direct inference of histories is computationally difficult on larger data sets, recent approaches such as ARGWeaver [44] are powerful enough to yield biological insights from some real data sets consisting of dozens or hundreds of genomes.

It may be efficient in some settings to use direct inference of a history to estimate the novelty profile. Indeed, once one has the history, computing the novelty profile is straightforward. In view of the strong assumptions underlying our persistent homology bounds on the novelty

profile, estimating the novelty profile via direct inference of the history seems to be an appealing alternative to the topological approach, for samples where it is computationally feasible. On the other hand, as emphasized in [8], one key advantage of the topological approach is computational efficiency for large genomic datasets.

## APPENDIX A. PROBABILITY THAT THE COALESCENT WITH RECOMBINATION GENERATES A GALLED TREE

**A.1. Overview of the Coalescent with Recombination.** The coalescent with recombination is a commonly used model of the evolutionary process generating a population genetic sample. For a detailed introduction to the coalescent with recombination, see [50]. Here, we give only a brief, informal description.

Instead of tracking an entire population, which may include millions or billions of reproducing organisms, a coalescent model tracks only the sampled individuals and their direct ancestors, up to their most recent common ancestor. We can think of the coalescent with recombination as a dynamical model that generates a phylogenetic graph, together with a time function on it, by proceeding backward in time. We start with an initial set of  $n$  vertices at some fixed final time, corresponding to  $n$  distinct lineages. As we proceed backwards in time, we can merge two lineages by adding a vertex of in-degree one and out-degree two, representing a common ancestor. We can also split a lineage into two distinct ones by adding a vertex of out-degree one and in-degree two, representing a recombinant. We require that the phylogenetic graph we create is rooted, so any split must eventually resolve itself by a merge further back in time. Once all lineages merge into the common ancestor, the graph-generating process stops; one can then generate the mutations at each vertex, using, e.g., a Poisson-type model as in Example 8.6. However, in this section, we will be concerned only with the underlying phylogenetic graph.

Two parameters are needed to specify the coalescent with recombination’s graph generation process: the number of leaves  $n$  and a recombination rate parameter  $\rho$ . The rate parameter equals twice the expected number of recombination events occurring in the entire population, per generation.

### A.2. Probability of Generating a Galled Tree as the Solution to a Linear System.

Let  $P(n, \rho)$  denote the probability that the coalescent with recombination generates a galled tree, given the parameters  $n$  and  $\rho$ . For fixed  $n$ , we derive a system of  $\mathcal{O}(n^2)$  linear equations, depending on  $\rho$ , whose solution gives an analytic expression for  $P(n, \rho)$  as a function of  $\rho$ . As  $n$  grows large, this expression becomes very complicated. But each linear system is sparse, so it is easy to solve for  $P(n, \rho)$  numerically for fixed values of  $n$  and  $\rho$ , provided  $n$  is not too large; see Fig. 17.

In the coalescent model, there are two types of *disallowed interactions* whose occurrence prevents the resulting graph from being a galled tree. First, after a split occurs, one of the two resulting branches may again split, prior to the resolution of the first split (Fig. 16a). Second, two unrelated splits may occur (resulting in four parental branches), after which a branch from one split joins with a branch from the other split (Fig. 16b).

To compute the probability  $P(n, \rho)$  that the coalescent generates a galled tree, we track the number of lineages ( $k$ ) and unresolved splits ( $s$ ) at each step of the process. The evolution

of  $k$  and  $s$  is described by a discrete-time Markov chain whose state space is the finite set  $T \cup \{X\}$ , where

$$T := \left\{ (k, s) \mid 1 \leq k \leq 2n, \max(0, k - n) \leq s \leq \frac{k}{2} \right\},$$

and  $X$  is an absorbing “failure state” that we enter into when a disallowed interaction occurs.  $P(n, \rho)$  is the probability that we eventually reach state  $(1, 0)$  in this Markov Chain, starting from  $(n, 0)$ .

To complete the description of the Markov chain, we first specify the transitions that can occur, and then specify the probabilities of each these. No self-transitions occur, with the exception that  $X$  and  $(1, 0)$  are both absorbing states. No two split or join events occur simultaneously in the coalescent; each transition thus corresponds to a single split or merge.

There are two types of splits: We may have a disallowed split, as described above and illustrated in (Fig. 16a); or the split may be allowed, in which case the state transition is  $(k, s) \mapsto (k + 1, s + 1)$  (Fig. 16c). There are three types of joins: We may have a disallowed join, as described above and illustrated in (Fig. 16b); a split may resolve itself (transition  $(k, s) \mapsto (k - 1, s - 1)$ ) (Fig. 16d), or two branches may join without altering any split (transition  $(k, s) \mapsto (k - 1, s)$ ) (Fig. 16e).

The transition probabilities are obtained as ratios of rates: In the coalescent, for  $k \geq 2$  the rate  $r_{\text{sp}}$  at which a split occurs is defined to be  $\rho k/2$ , and the rate  $r_{\text{jn}}$  at which two branches join together is defined to be  $k(k - 1)/2$  [50]. The total rate of a split or merge is then  $r_{\text{tot}} := r_{\text{sp}} + r_{\text{jn}}$ , i.e.,  $r_{\text{tot}} = k(k + \rho - 1)/2$ .

The rate  $r_{\text{asp}}$  of an allowed split is the product of  $r_{\text{sp}}$  with the fraction of branches for which splits are allowed, i.e.,  $r_{\text{asp}} = \rho(k/2 - s)$ . For  $k \geq 2$ , let  $p_S(k, s)$  denote the probability of an allowed split (i.e., a transition  $(k, s) \mapsto (k, s + 1)$ ). Then

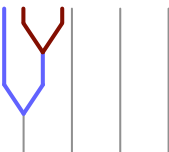
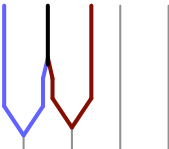
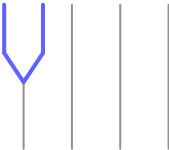
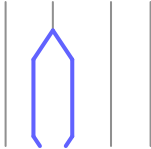
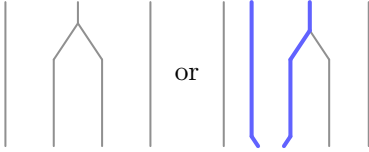
$$p_S(k, s) = r_{\text{asp}}/r_{\text{tot}} = \frac{\rho(k - 2s)}{k(k + \rho - 1)}.$$

The other transition probabilities are obtained analogously. We denote them as follows:  $p_R(k, s)$  is the probability that a split is resolved;  $p_J(k, s)$  is the probability of an allowed join that does not resolve a split; and  $p_X(k, s)$  is the probability of a disallowed interaction of either type. The formulas for these are given in Fig. 16.

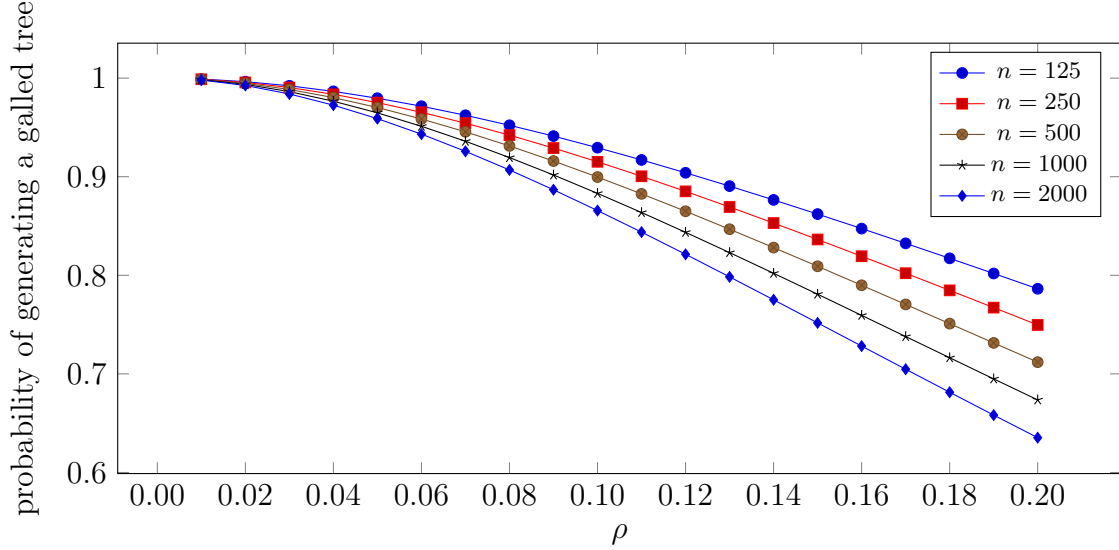
For  $(k, s) \in T$ , let  $f(k, s)$  be the probability that we eventually generate a galled tree, given that the current state is  $(k, s)$ , and for  $(k, s) \in \mathbb{Z}^2 \setminus T$ , let  $f(k, s) = 0$ . Thus,  $P(n, \rho) = f(n, 0)$ . The  $f(k, s)$  satisfy the linear system

$$\begin{aligned} f(k, s) &= p_S(k, s)f(k + 1, s + 1) \\ &\quad + p_R(k, s)f(k - 1, s - 1) \\ &\quad + p_J(k, s)f(k - 1, s) \quad \text{for } (k, s) \in T \setminus \{(1, 0)\}, \\ f(1, 0) &= 1. \end{aligned} \tag{3}$$

Fig. 17 gives values of  $P(n, \rho)$  for several choices of  $\rho$  and  $n$ , obtained by solving this linear system numerically. We observe that for fixed  $n$ ,  $P(n, \rho)$  tends to 1 as  $\rho$  tends to 0. Varying both parameters,  $1 - P(n, \rho)$  appears to decrease on the order  $\rho^2 (\log n)^2$  for large  $n$  and small  $\rho$ .

	<u>Probability</u>	<u>Result</u>	
<b>a</b>		$\frac{2\rho s}{k(k+\rho-1)} + \frac{4s(s-1)}{k(k+\rho-1)} \quad p_X(k, s)$	Process terminates (genealogy is not a galled tree.)
<b>b</b>			
<b>c</b>		$p_S(k, s) = \frac{\rho(k-2s)}{k(k+\rho-1)}$	$(k, s) \mapsto (k+1, s+1)$
<b>d</b>		$p_R(k, s) = \frac{2s}{k(k+\rho-1)}$	$(k, s) \mapsto (k-1, s-1)$
<b>e</b>		$p_J(k, s) = \frac{k(k-1) - 2s(2s-1)}{k(k+\rho-1)}$	$(k, s) \mapsto (k-1, s)$

**Figure 16.** Possible events in the coalescent process described in the text. Light gray lines indicate ordinary lineages, medium blue and dark red lines indicate pairs of parental lineages that split from their recombinant child lineages, and the black line indicates a join between the red and blue lineages. Each diagram is read upwards (back in time). To the right of each diagram are the corresponding transition probability and the resulting change in state of the Markov chain, in terms of population-scaled recombination rate  $\rho$ , number of lineages  $k$ , and number of unresolved splits  $s$ . Where multiple events are depicted, the probability is given only for the topmost (most ancient) event. (a) A disallowed split occurs when a lineage that has already split (light blue) splits again (dark red). (b) A disallowed join occurs when parental lineages from two separate splitting events (blue, red) join together (black). (c) An allowed split. (d) A join that resolves a split. (e) An allowed join that does not resolve a split.



**Figure 17.** Probability that the coalescent with recombination yields a galled tree, for several values of the recombination rate parameter  $\rho$  and the number of sampled genomes  $n$ . We see that for fixed  $n$ , the probability of obtaining a galled tree tends to 1 as  $\rho$  tends to 0.

To give a sense of scale for human population genetics (species effective population size  $N_e \approx 10^4$ , recombination rate  $c \approx 10^{-8} \text{ bp}^{-1}$ ), in a sample of 125 individuals ( $n = 250$  haploid genomes), to ensure that disallowed interactions occur with probability less than 0.1 ( $P(n, \rho) > 0.9$ ), the requirement shown in Fig. 17 is  $\rho < 0.11$ . Using  $\rho = 4N_e c L$ , where  $L$  is the length of the genome segment analyzed, this requirement becomes  $L < 275 \text{ bp}$ . In a sample of one thousand ( $n = 2000$  in Fig. 17), the requirement is instead  $\rho < 0.08$ , or  $L < 200 \text{ bp}$ . Extrapolating from the apparent  $\rho^2 (\log n)^2$  scaling, in a sample of one million ( $n = 2 \times 10^6$ ), the requirement is  $\rho < 0.04$ , or  $L < 100 \text{ bp}$ .

## APPENDIX B. SUBSAMPLES RARELY VIOLATE OUR THEORETICAL BOUNDS FOR COMPLETE SAMPLES

Example 6.19 makes clear that Theorem 6.18 (i), our persistent homology lower bound on the novelty profile of an evolutionary history, does not hold for arbitrary samples  $S$  of the history  $\mathcal{E}$ . Nevertheless, one might hope that violations of Theorem 6.18 are relatively rare. Here we use simulations to explore this question in the coalescent model with recombination.

Our computations focus on how often the number of intervals in  $\mathcal{B}_1(S)$  exceeds the number of recombinants in  $\mathcal{E}$ . We find that in our simulations, this happens quite rarely, though it does occur. We did not consider the frequency of other kinds of violations of Theorem 6.18 for subsamples, though it would be interesting to do so.

We simulated over 42,000 evolutionary histories, assuming a constant population size and using parameters  $n = 10, 15$ , or  $20$ ;  $\rho = 1, 2, 3$ , or  $4$ ; and  $\theta = 5, 10$ , or  $30$ , where  $\theta$  is the mutation rate parameter for the coalescent [50]. We used rejection sampling, retaining only those histories that were indexed by galled trees. Consistent with Section 2, we used



an infinite-sites model of mutation. Each genetic site in a recombinant offspring inherited the state of a parent with probability one-half. For each simulated history, we counted the number of detectable recombination events that took place – defined as the number of events giving rise to a recombinant that generates an incompatibility according to the four-gamete test [32]. We then computed the maximum  $|\mathcal{B}_1(S)|$  among 2500 random subsamples  $S$  of the history, with 5 to 30 genomes in each subsample. Among all simulations, nine had a maximum  $|\mathcal{B}_1(S)|$  greater than the true number of detectable recombination events (Table 1). The rarity of this violation suggests that counterexamples such as Example 6.19 may be uncommon in actual population-genetic data.

**Table 1.** Counts of coalescent simulations, by number of detectable recombination events and maximum number of intervals in  $\mathcal{B}_1(S)$  among all subsamples  $S$ . Gray shading indicates “false positive” scenarios in which more events are detected than actually occurred. Nine of the simulations exhibited a false positive (shown in red).

		Maximum $ \mathcal{B}_1(S) $ among all subsamples $S$					
		0	1	2	3	4	5
Number of detectable recombination events in ARG	0	17953	0	0	0	0	0
	1	9346	10049	9	0	0	0
	2	1159	2632	714	0	0	0
	3	126	395	182	14	0	0
	4	16	48	39	3	0	0
	5	0	5	2	0	0	0

## REFERENCES

- [1] M. Adamaszek and H. Adams. The vietoris–rips complexes of a circle. *Pacific Journal of Mathematics*, 290(1):1–40, 2017.
- [2] M. Adamaszek, H. Adams, E. Gasparovic, M. Gommel, E. Purvine, R. Sazdanovic, B. Wang, Y. Wang, and L. Ziegelmeier. Vietoris-rips and cech complexes of metric gluings. *arXiv preprint arXiv:1712.06224*, 2017. To appear in SoCG 2018.
- [3] M. Adamaszek, H. Adams, and S. Reddy. On vietoris–rips complexes of ellipses. *Journal of Topology and Analysis*, pages 1–30, 2017.
- [4] M. Arenas, G. Valiente, and D. Posada. Characterization of reticulate networks based on the coalescent with recombination. *Molecular biology and evolution*, 25(12):2517–2520, Dec. 2008.
- [5] U. Bauer and M. Lesnick. Induced matchings and the algebraic stability of persistence barcodes. *Journal of Computational Geometry*, 6(2):162–191, 2015.
- [6] D. Bezemer, A. van Sighem, V. V. Lukashov, L. van der Hoek, N. Back, R. Schuurman, C. A. B. Boucher, E. C. J. Claas, M. C. Boerlijst, R. A. Coutinho, F. de Wolf, and ATHENA observational cohort. Transmission networks of HIV-1 among men having sex with men in the Netherlands. *AIDS (London, England)*, 24(2):271–282, Jan. 2010.
- [7] A. J. Blumberg and M. Lesnick. Universality of the homotopy interleaving distance. *arXiv preprint arXiv:1705.01690*, 2017.
- [8] P. G. Cámara, A. J. Levine, and R. Rabadán. Inference of ancestral recombination graphs through topological data analysis. *PLoS Comput Biol*, 12(8):e1005071, 2016.

- [9] P. G. Cámara, D. I. S. Rosenbloom, K. J. Emmett, A. J. Levine, and R. Rabadán. Topological Data Analysis Generates High-Resolution, Genome-wide Maps of Human Recombination. *Cell systems*, 3(1):83–94, July 2016.
- [10] G. Carlsson. Topology and data. *American Mathematical Society*, 46(2):255–308, 2009.
- [11] G. Carlsson. Topological pattern recognition for point cloud data. *Acta Numerica*, 23:289–368, May 2014.
- [12] J. Chan, G. Carlsson, and R. Rabadán. Topology of viral evolution. *Proceedings of the National Academy of Sciences*, 110(46), November 2013.
- [13] F. Chazal, D. Cohen-Steiner, M. Glisse, L. Guibas, and S. Oudot. Proximity of persistence modules and their diagrams. In *Proceedings of the 25<sup>th</sup> annual symposium on Computational geometry*, pages 237–246. ACM, 2009.
- [14] F. Chazal, D. Cohen-Steiner, L. Guibas, F. Méholi, and S. Oudot. Gromov-Hausdorff stable signatures for shapes using persistence. In *Proceedings of the Symposium on Geometry Processing*, pages 1393–1403. Eurographics Association, 2009.
- [15] F. Chazal, V. de Silva, M. Glisse, and S. Oudot. *The Structure and Stability of Persistence Modules*. Springer International Publishing, 2016.
- [16] F. Chazal, V. De Silva, and S. Oudot. Persistence stability for geometric complexes. *Geometriae Dedicata*, 173(1):193–214, 2014.
- [17] D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Discrete and Computational Geometry*, 37(1):103–120, 2007.
- [18] T. F. Cooper. Recombination speeds up adaptation by reducing competition between beneficial mutations in populations of *Escherichia coli*. *PLoS Biology*, 5(9):1899–1905, 2007.
- [19] W. Crawley-Boevey. Decomposition of pointwise finite-dimensional persistence modules. *Journal of Algebra and Its Applications*, 14(05):1550066, 2015.
- [20] J. Davies and D. Davies. Origins and evolution of antibiotic resistance. *Microbiology and Molecular Biology Reviews*, 74(3):417–433, 2010.
- [21] R. Durrett. *Probability: theory and examples*. Cambridge university press, 2010.
- [22] H. Edelsbrunner and J. Harer. *Computational topology: an introduction*. American Mathematical Society, 2010.
- [23] K. Emmett, D. Rosenbloom, P. Cámara, and R. Rabadán. Parametric inference using persistence diagrams: A case study in population genetics. *Proc. 31st Intl. Conf. Machine Learning*, 2014.
- [24] K. J. Emmett and R. Rabadán. Characterizing Scales of Genetic Recombination and Antibiotic Resistance in Pathogenic Bacteria Using Topological Data Analysis. *Lecture Notes in Computer Science*, 8609:540–551, 2014.
- [25] R. Forman. Morse theory for cell complexes. *Advances in mathematics*, 134(1):90–145, 1998.
- [26] R. Forman. A user’s guide to discrete Morse theory. *Sém. Lothar. Combin.*, 48:Art. B48c, 35, 2002.
- [27] P. J. Gerrish, A. Colato, A. S. Perelson, and P. D. Sniegowski. Complete genetic linkage can subvert natural selection. *Proc. Nat’l Acad. Sci. USA*, 104(15):6266–6271, 2007.
- [28] D. Gusfield. *ReCombinatorics: the algorithmics of ancestral recombination graphs and explicit phylogenetic networks*. MIT Press, 2014.
- [29] D. Gusfield, S. Eddhu, and C. Langley. Efficient reconstruction of phylogenetic networks with constrained recombination. In *Bioinformatics Conference, 2003. CSB 2003. Proceedings of the 2003 IEEE*, pages 363–374. IEEE, 2003.
- [30] S. Harker, M. Kramar, R. Levanger, and K. Mischaikow. A comparison framework for interleaved persistence modules. *arXiv preprint arXiv:1801.06725*, 2018.
- [31] A. Hatcher. *Algebraic topology*. Cambridge University Press, 2002.

- [32] R. R. Hudson and N. L. Kaplan. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*, 111(1):147–164, Sept. 1985.
- [33] D. H. Huson, R. Rupp, and C. Scornavacca. *Phylogenetic networks: concepts, algorithms and applications*. Cambridge University Press, 2010.
- [34] T. Ito, J. N. S. S. Couceiro, S. Kelm, L. G. Baum, S. Krauss, M. R. Castrucci, I. Donatelli, H. Kida, J. C. Paulson, R. G. Webster, and Y. Kawaoka. Molecular basis for the generation in pigs of influenza a viruses with pandemic potential. *Journal of Virology*, 72(9):7367–7373, 1998.
- [35] M. Kahle. Random geometric complexes. *Discrete & Computational Geometry*, 45(3):553–573, 2011.
- [36] O. Kallenberg. *Foundations of modern probability*. Springer Science & Business Media, 2006.
- [37] D. Kozlov. *Combinatorial algebraic topology*, volume 21 of *Algorithms and Computation in Mathematics*. Springer, Berlin, 2008.
- [38] S. L. Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- [39] M. J. McDonald, D. P. Rice, and M. M. Desai. Sex speeds adaptation by altering the dynamics of molecular evolution. *Nature*, 531:233–236, 2016.
- [40] J. R. Munkres. *Elements of algebraic topology*, volume 2. Addison-Wesley Menlo Park, 1984.
- [41] S. R. Myers and R. C. Griffiths. Bounds on the minimum number of recombination events in a sample history. *Genetics*, 163(1):375–394, Jan. 2003.
- [42] T. Nora, C. Charpentier, O. Tenaillon, C. Hoede, F. Clavel, and A. J. Hance. Contribution of recombination to the evolution of human immunodeficiency viruses expressing resistance to antiretroviral treatment. *Journal of virology*, 81(14):7620–7628, 2007.
- [43] S. Y. Oudot. *Persistence Theory: From Quiver Representations to Data Analysis*. Number 209 in AMS Mathematical Surveys and Monographs. American Mathematical Society, 2015.
- [44] M. D. Rasmussen, M. J. Hubisz, I. Gronau, and A. Siepel. Genome-Wide Inference of Ancestral Recombination Graphs. *PLoS Genetics*, 10(5):e1004342, May 2014.
- [45] A. Solovyov, G. Palacios, T. Brieze, W. I. Lipkin, and R. Rabadan. Cluster analysis of the origins of the new influenza a (h1n1) virus. *Euro surveillance: bulletin European sur les maladies transmissibles = European communicable disease bulletin*, 14(21), 2009.
- [46] Y. S. Song, Y. Wu, and D. Gusfield. Efficient computation of close lower and upper bounds on the minimum number of recombinations in biological sequence evolution. *Bioinformatics (Oxford, England)*, 21 Suppl 1(Suppl 1):i413–22, June 2005.
- [47] T. Stadler, R. Kouyos, V. von Wyl, S. Yerly, J. Boni, P. Burgisser, T. Klimkait, B. Joos, P. Rieder, D. Xie, H. F. Gunthard, A. J. Drummond, S. Bonhoeffer, and the Swiss HIV Cohort Study. Estimating the Basic Reproductive Number from Viral Sequence Data. *Molecular biology and evolution*, 29(1):347–357, Dec. 2011.
- [48] The Swiss HIV Cohort Study. Cohort profile: the Swiss HIV Cohort study. *International Journal of Epidemiology*, 39:1179–1189, 2010.
- [49] V. Trifonov, H. Khiabani, R. Rabadan, et al. Geographic dependence, surveillance, and origins of the 2009 influenza a (h1n1) virus. *New England Journal of Medicine*, 361(2):115–119, 2009.
- [50] J. Wakeley. *Coalescent Theory*. An Introduction. Roberts & Co., 2007.
- [51] L. Wang, K. Zhang, and L. Zhang. Perfect phylogenetic networks with recombination. *Journal of Computational Biology*, 8(1):69–78, 2001.

LESNICK: PRINCETON UNIVERSITY, PRINCETON, NJ, USA  
*E-mail address:* mlesnick@princeton.edu

RABADÁN: COLUMBIA UNIVERSITY, NEW YORK, NY, USA  
*E-mail address:* rr2579@cumc.columbia.edu

ROSENBLOOM: MERCK RESEARCH LABORATORIES, RAHWAY, NJ, USA  
*E-mail address:* daniel.rosenbloom@merck.com