

# **Topology of Reticulate Evolution**

**Kevin Joseph Emmett**

Submitted in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy  
in the Graduate School of Arts and Sciences

**COLUMBIA UNIVERSITY**

2016

© 2016

Kevin Joseph Emmett

All Rights Reserved

# ABSTRACT

## Topology of Reticulate Evolution

**Kevin Joseph Emmett**

The standard representation of evolutionary relationships is a bifurcating tree. However, many types of genetic exchange, collectively referred to as *reticulate evolution*, involve processes that cannot be modeled as trees. Increasing genomic data has pointed to the prevalence of reticulate processes, particularly in microorganisms, and underscored the need for new approaches to capture and represent the scale and frequency of these events.

This thesis contains results from applying new techniques from applied and computational topology, under the heading *topological data analysis*, to the problem of characterizing reticulate evolution in molecular sequence data. First, we develop approaches for analyzing sequence data using topology. We propose new topological constructions specific to molecular sequence data that generalize standard constructions such as Vietoris-Rips. We draw on previous work in phylogenetic networks and use homology to provide a quantitative measure of reticulate events. We develop methods for performing statistical inference using topological summary statistics.

Next, we apply our approach to several types of molecular sequence data. First, we examine the mosaic genome structure in phages. We recover inconsistencies in existing morphology-based taxonomies, use a network approach to construct a genome-based representation of phage relationships, and identify conserved gene families within phage popu-

lations. Second, we study influenza, a common human pathogen. We capture widespread patterns of reassortment, including nonrandom cosegregation of segments and barriers to subtype mixing. In contrast to traditional influenza studies, which focus on the phylogenetic branching patterns of only the two surface-marker proteins, we use whole-genome data to represent influenza molecular relationships. Using this representation, we identify unexpected relationships between divergent influenza subtypes. Finally, we examine a set of pathogenic bacteria. We use two sources of data to measure rates of reticulation in both the core genome and the mobile genome across a range of species. Network approaches are used to represent the population of *S. aureus* and analyze the spread of antibiotic resistance genes. The presence of antibiotic resistance genes in the human microbiome is investigated.

# Contents

<b>Contents</b>	<b>i</b>
<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Molecular Evolution and the Tree Paradigm . . . . .	2
1.2 Reticulate Processes and the Universal Tree . . . . .	6
1.3 Evolution as a Topological Space . . . . .	8
1.4 Thesis Organization . . . . .	11
<b>2 Background</b>	<b>13</b>
2.1 Evolutionary Biology and Genomics . . . . .	13
2.1.1 Genes and Genomes . . . . .	14
2.1.2 Evolutionary Processes . . . . .	14
2.1.3 Mathematical Models of Evolution . . . . .	19
2.1.4 Phylogenetic Methods . . . . .	22
2.2 Topological Data Analysis . . . . .	29
2.2.1 Preliminaries . . . . .	33
2.2.2 Persistent Homology . . . . .	38
2.2.3 Mapper . . . . .	48
2.3 Applying TDA to Molecular Sequence Data . . . . .	50
2.3.1 Topology of Tree-like Metrics . . . . .	51
2.3.2 The Fundamental Unit of Reticulation . . . . .	51
2.3.3 A Complete Example . . . . .	53
2.3.4 The Space of Trees, Revisited . . . . .	54
<b>I Theory</b>	<b>57</b>
<b>3 Quantifying Reticulation Using Topological Complex Constructions Beyond Vietoris-Rips</b>	<b>59</b>
3.1 Introduction . . . . .	59
3.2 Sensitivity of the Vietoris-Rips Construction . . . . .	60
3.3 The Median Complex Construction . . . . .	62

3.4	Čech Complex Construction as an Optimization Problem . . . . .	67
3.5	Conclusions . . . . .	70
<b>4</b>	<b>Parametric Inference using Persistence Diagrams</b>	<b>73</b>
4.1	Introduction . . . . .	73
4.2	Warmup: Gaussian Random Fields . . . . .	75
4.3	The Coalescent Process . . . . .	75
4.4	Statistical Model . . . . .	77
4.5	Coalescent Simulations . . . . .	80
4.6	Conclusions . . . . .	81
<b>II</b>	<b>Applications</b>	<b>83</b>
<b>5</b>	<b>Phage Mosaicism</b>	<b>85</b>
5.1	Introduction . . . . .	85
5.2	Data . . . . .	88
5.3	Measuring Phage Mosaicism with Persistent Homology . . . . .	89
5.4	Representing Phage Relationships with Mapper . . . . .	91
5.5	Conclusions . . . . .	96
<b>6</b>	<b>Reassortment in Influenza Evolution</b>	<b>101</b>
6.1	Introduction . . . . .	101
6.2	Influenza Virology . . . . .	103
6.3	Influenza Reassortment . . . . .	103
6.4	Nonrandom Association of Genome Segments . . . . .	106
6.5	Multiscale Flu Reassortment . . . . .	110
6.6	Conclusions . . . . .	111
<b>7</b>	<b>Reticulate Evolution in Pathogenic Bacteria</b>	<b>113</b>
7.1	Introduction . . . . .	113
7.2	Evolutionary Scales of Recombination in the Core Genome . . . . .	115
7.3	Protein Families as a Proxy for Genome Wide Reticulation . . . . .	118
7.4	Antibiotic Resistance in <i>Staphylococcus aureus</i> . . . . .	121
7.5	Microbiome as a Reservoir of Antibiotic Resistance Genes . . . . .	122
7.6	Conclusions . . . . .	124
<b>8</b>	<b>Conclusions</b>	<b>127</b>
<b>Bibliography</b>		<b>129</b>

# List of Figures

1.1	Charles Darwin and the Evolutionary Tree . . . . .	3
1.2	Carl Woese’s Three Domain Tree of Life . . . . .	5
1.3	Ford Doolittle’s Reticulate Tree of Life . . . . .	8
1.4	Topological equivalence of the coffee mug and the donut . . . . .	9
1.5	Treelike and reticulate phylogenies . . . . .	11
2.1	Viral recombination and reassortment . . . . .	16
2.2	Three modes of bacterial reticulation . . . . .	17
2.3	Two models for simulating evolutionary data . . . . .	21
2.4	Rooted vs. Unrooted trees . . . . .	23
2.5	The four point condition for additivity . . . . .	25
2.6	Counting Tree Topologies . . . . .	27
2.7	Tree Space . . . . .	29
2.8	Example of a Splits Network . . . . .	30
2.9	Simplices: The building blocks of topological complexes . . . . .	34
2.10	Simplicial Complex: A discrete topological space . . . . .	34
2.11	Relationship between the chain group, cycle group, and boundary group . . . . .	36
2.12	Simplicial Homology . . . . .	36
2.13	Vietoris-Rips and ČechComplexes . . . . .	38
2.14	Multiscale Topological Structure . . . . .	41
2.15	Barcode Diagram for the Two Circles Example . . . . .	42
2.16	The Persistence Pipeline . . . . .	42
2.17	Level-Set Filtrations . . . . .	43
2.18	The stability of persistent homology under noise . . . . .	45
2.19	Confidence sets on the persistence diagram . . . . .	46
2.20	Persistence landscape of the persistence diagram . . . . .	47
2.21	Dimensionality Reduction for EDA . . . . .	48
2.22	The Mapper Algorithm . . . . .	49
2.23	Fundamental Unit of Reticulation . . . . .	53
2.24	Applying TDA to Molecular Sequence Data . . . . .	55
2.25	Topology expands tree space to include reticulation . . . . .	56
3.1	Two examples of reduced sensitivity of the Vietoris-Rips Complex . . . . .	62
3.2	The Median Operation on Binary Sequences . . . . .	63
3.3	The Median Complex Recovers Reticulation in Example One . . . . .	64

3.4	The Median Complex Recovers Reticulation in Example Two . . . . .	65
3.5	Recombination in <i>D. melanogaster</i> . . . . .	67
3.6	Species hybridization in genus <i>Ranunculus</i> . . . . .	68
4.1	Barcode Diagrams for Three Coalescent Simulations . . . . .	75
4.2	Gaussian Random Fields with Exponential Covariance Structure . . . . .	76
4.3	Gaussian Random Field Summary Statistics . . . . .	76
4.4	Distributions of statistics defined on the $H_1$ persistence diagram for different model parameters . . . . .	78
4.5	Inference of recombination rate $\rho$ using topological information . . . . .	81
4.6	Comparing traditional estimators of $\rho$ to . . . . .	82
5.1	Inconsistency of morphological classifications in bacteriophage . . . . .	87
5.2	Summary annotations of 306 bacteriophage strains used in this study . . . . .	89
5.3	S306 Bacteriophage Barcode Diagram . . . . .	91
5.4	Caudovirales $H_0$ dendrogram . . . . .	92
5.5	Caudovirales Barcode Diagrams . . . . .	93
5.6	Phage Mapper Network . . . . .	94
5.7	Phage Network Colored by Taxonomic Family . . . . .	95
5.8	Modularity Scores for Different Divisions of the Phage Network . . . . .	96
5.9	Phage Network Colored by Host . . . . .	97
5.10	Phage Network with MCL Clustering . . . . .	99
6.1	Structure of an influenza virus particle . . . . .	104
6.2	Influenza Dataset Statistics . . . . .	105
6.3	Influenza Genome Segment Barcodes . . . . .	107
6.4	Influenza Concatenated Genome Barcode . . . . .	108
6.5	Influenza Networks By HA Subtype . . . . .	109
6.6	Influenza Nonrandom Reassortment . . . . .	110
6.7	$H_1$ persistence diagram computed from an avian influenza dataset. . . . .	112
7.1	Core genome exchange in <i>K. pneumoniae</i> and <i>S. enterica</i> . . . . .	117
7.2	$H_1$ persistence diagram for twelve pathogenic strains using MLST profile data . . . . .	118
7.3	Core genome reticulation patterns in pathogenic bacteria from MLST profiles . . . . .	119
7.4	Genome-wide reticulation patterns in pathogenic bacteria from protein annotations . . . . .	120
7.5	FigFam similarity network of <i>S. aureus</i> . . . . .	122
7.6	FigFam similarity network of the gastrointestinal tract . . . . .	124

# List of Tables

2.1	Reticulate processes in biology across kingdoms . . . . .	18
2.2	Dictionary connecting algebraic topology and evolutionary biology . . . . .	54
3.1	Čech Homology of Hypercube . . . . .	70
5.1	Phage families defined by the ICTV . . . . .	86
5.2	Phage Network MCL clustering annotations and representative protein families . .	98
6.1	Influenza Protein Segments . . . . .	104
7.1	List of pathogenic bacteria selected for study . . . . .	115



# Chapter 1

## Introduction

Charles Darwin's *On the Origin of Species* contains a single figure, depicting the ancestry of species as a branching genealogical tree, or *phylogeny* [41] (see Figure 1.1). Darwin argued that evolution was mediated by descent with modification; that is, the gradual change in heritable traits under the pressure of natural selection. Since that time, the tree structure has been the dominant framework to understand, visualize, and communicate discoveries about evolution. Indeed, an important aim of evolutionary biology has been expanding the *universal tree of life*, the set of evolutionary relationships among all extant and extinct organisms on Earth [19].

Traditionally, evolutionary relationships were established on the basis of phenotype, i.e. the observable traits of each organism. With the advent of molecular models of evolution and rapidly increasing genomic sequence data, the genotype has supplanted phenotype as the primary focus of evolutionary studies. Molecular phylogenetics has become established as the standard tool for inferring phylogenetic relationships. However, a phylogenetic tree is accurate only if the Darwinian model of descent with modification is the sole process driving evolution. It has long been recognized that there exist alternative evolutionary processes that can allow organisms to directly exchange genetic material [5]. Notable examples include horizontal gene transfer in bacteria [123], species hybridization in plants [4], and meiotic

recombination in eukaryotes [37]. Collectively, these processes are referred to as *reticulate evolution*. Reticulate evolution stand in contrast to the paradigm of tree-like diversification, an example of *clonal evolution*.<sup>1</sup> Increasing genomic data, powered by new high-throughput sequencing technologies, has shown that these reticulate processes are more prevalent than originally expected [18]. For some, this has called into question the tree of life hypothesis as an organizing principle and prompted the search for new ways of representing evolutionary relationships [45, 124, 94].

This thesis presents a new approach to quantifying and representing reticulate evolutionary processes using recently developed ideas from algebraic and computational topology. The methods we employ fall under the collective heading of *topological data analysis* (henceforth TDA), a new branch of applied topology concerned with inferring structure in high-dimensional data [28]. The thesis consists of three aims: (1) introduce the methods of TDA and their application to biological and genomic data; (2) develop approaches tailored to the unique features of molecular sequence data; and (3) apply these approaches to a range of biological problems in which reticulate processes are believed to play an important role.

In the following brief introduction, we survey salient aspects of molecular evolution, the tree paradigm, and the challenges posed by reticulate processes. We then introduce the idea of representing evolution as a topological space and give a flavor of the results to be discussed.

## 1.1 Molecular Evolution and the Tree Paradigm

The combination of Darwin's theory of natural selection with Mendelian genetics led to the *modern evolutionary synthesis*, outlined in the first half of the twentieth century in pioneering works by Ronald Fisher, Sewall Wright, JBS Haldane, and others.<sup>2</sup> The modern

---

<sup>1</sup>Clonal and reticulate evolution are also known by the terms *vertical* and *horizontal* evolution, respectively.

<sup>2</sup>See [84] and [73] for comprehensive historical reviews.

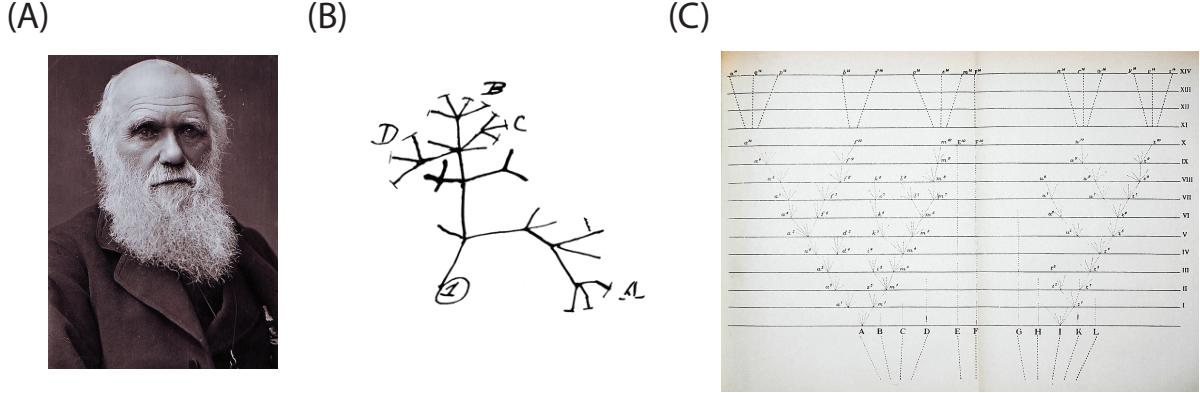


Figure 1.1: (A) Charles Darwin in 1881. Photograph by Herbert Rose Barraud. (B) An sketch from Darwin’s notebooks (circa 1837) showing an early conception of the evoutionary tree. (C) The only figure in Darwin’s *On the Origin of Species*. In *The Origin*, Darwin argued for descent with modification and natural selection as the driving processes underscoring evolution. In this figure, Darwin illustarted his idea for how diverging species would result in a tree structure. Reproduced from [41].

synthesis was based largely on an analysis of distributions of allele frequencies in distinct populations, the purview of classical population genetics. The field was placed on a molecular foundation with Watson and Crick’s discovery of the DNA double-helix in 1953 [150]. These developments led to the establishment of *molecular evolution*, the analysis of how processes such as mutation, drift, and recombination act to induce changes in populations and species.

The information underlying an organism’s form and function is encoded in its genome, the complete sequence of DNA (or RNA) contained in each cell. The genome can be represented as a string of nucleotides, indexed by position. Embedded within the genome are regions defining the genes which code for functional proteins, as well as non-coding regions which have as-yet unknown function.<sup>3</sup> When an organism reproduces, either sexually or asexually, a complete copy of this genomic information is passed to the offspring. Because the molecular mechanisms that control this copying are not exact, errors in replication are introduced. These errors can take the form of single point mutations (or single nucleotide polymorphisms,

---

<sup>3</sup>In humans, only 1.5% of the genome is protein-coding, the rest largely non-functional [100]. Up to 5-8% of the human genome is believed to consist of endogenous retroviruses, dead viruses which have integrated their genome into the human genome [14].

SNPs), small insertions and deletions of a few nucleotides (indels), or larger effects including copy number variations (CNVs) and chromosomal duplications.<sup>4</sup> Under the neutral theory of evolution, the majority of these errors will have very little impact, either positive or negative, on the descendant organism. A small fraction of mutations will result in an appreciable fitness difference compared to other organisms, and it is on these organisms that natural selection will act.

While molecular biology has largely focused on the biochemical and biophysical mechanisms underlying these processes, *molecular phylogenetics* has focused on the comparative analysis of macromolecular sequences to infer genealogical and evolutionary relationships. Molecular phylogenetics began with Emile Zuckerkandl and Linus Pauling’s recognition in the early 1960’s that the information encoded in a set of molecular sequences could itself be used as a document of evolutionary history [162, 163]. It became clear that given two sequenced organisms, counting the differences between their respective sequences could be used as a quantitative measure of the amount of evolutionary divergence between the two organisms. If one has a larger set of sequenced organisms, computing the complete set of pairwise distances yields a *distance matrix*. From the distance matrix, one can then attempt to associate a tree to the data such that pairwise distances along the tree are close to the measured pairwise distances from the sequences. Walter Fitch and Emanuel Margoliash popularized this approach by constructing a weighted least squares approach to fitting phylogenetic trees from distances [65]. Since that time, the development of numerical approaches for inferring evolutionary relationships has evolved into a mature discipline and the use of molecular sequence data to infer phylogeny has become a standard practice across a wide range of biology and ecology. While other approaches to tree inference have been developed, including parsimony, maximum likelihood (ML), and Bayesian methods, we will focus on distance matrix methods because of their close relationship to the topological ideas

---

<sup>4</sup>Mutation rate vary across species: in humans,  $10^{-8}$  per site per generation [116]; in bacteria and unicellular eukaryotes, between  $10^{-9}$  and  $10^{-10}$  per site per generation; in DNA viruses, between  $10^{-6}$  and  $10^{-8}$  per site per generation [48].

# Phylogenetic Tree of Life

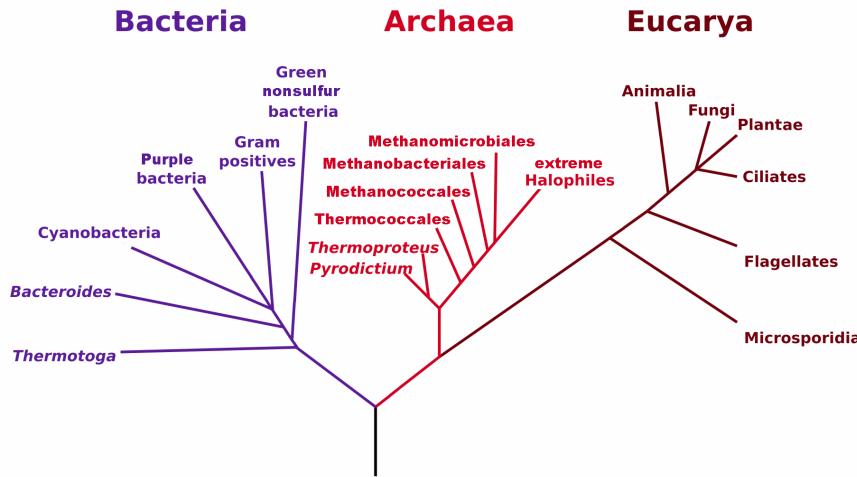


Figure 1.2: Carl Woese’s three domain tree of life. Using 16S subunit ribosomal RNA, Carl Woese identified archaea as a distinct phylogenetic domain. Previously, based on morphological similarity (specifically, unicellular and lacking a nucleus), archaea had been grouped with bacteria. This result was an early success for molecular phylogenetics and the use of conserved gene segments for molecular classification. Figure adapted from [156].

we employ.

One important early result from molecular phylogenetics was Carl Woese’s organization of bacteria, eukarya, and archaea into the three domains of life [155]. Prior to Woese, there were two recognized domains of life: prokaryotes, single-celled organisms lacking a nucleus, and eukaryotes, multi-celled organisms with an enveloped nucleus. Using 16S subunit ribosomal RNA sequencing, Woese discovered that the prokaryotic domain actually split into two evolutionarily distinct groups. One of these, which he termed *archaeabacteria* was more closely related to eukaryotes than were the rest of the prokaryotes. This led to the three-domain system of life (Figure 1.2).

This work had several important consequences. First, it established the use of molecular data to inform about large-scale patterns of evolutionary history. Using only morphological data had led to an inconsistent classification of archaea. Second, it positioned 16S rRNA profiling as the primary source of data for use in comparative genomics. The use of this

genomic region was justified on the basis of being one of the few universal gene segments that is conserved across all species. Constructing a universal tree is predicated on there being orthologous genes, i.e. shared genes related through speciation events, that can provide a common foundation for comparative study. Finally, it solidified the tree paradigm as an organizing principle for relating extant species. Even though reticulate processes had been known since the early twentieth century<sup>5</sup>, the idea that evolutionary relationships should be described by a bifurcating tree had been paramount since Darwin. Reticulate processes were either ignored completely, or expected to occur at such low frequencies that they need not be considered.

## 1.2 Reticulate Processes and the Universal Tree

Despite the significant impact of Woese’s observation, there remained a subtle difficulty, which Woese himself would come to contemplate in later work [154, 72]. Woese’s phylogeny was based on only 1,500 nucleotides in the ribosomal RNA, less than 1% of the total length of a typical bacterial genome (see [40]). Even more striking, this accounts for less than 0.00005% of the human genome. While recent work has developed approaches for constructing reference trees from larger gene sets [36], the fact remains that the vast majority of genomic information is *not* incorporated into the tree.

The reason for this situation is twofold. First, not all genes are shared universally across all species. In constructing a phylogenetic tree using sequence data, only genes that are present across all species are informative. Second, even among universal genes, the presence of reticulate evolutionary processes will confound systematic analysis. The model of a bifurcating tree will be consistent only if all loci share the same pattern of bifurcation. When organisms can exchange genetic material by means other than direct reproduction, the ancestral relationships between organisms will depend on which genomic regions are

---

<sup>5</sup>Beginning with Frederick Griffith’s experiments in 1928 showing that non-virulent strains of *Streptococcus pneumoniae* could acquire virulence factors by being exposed to dead virulent strains.

used. If two different genomic regions were analyzed, two different tree topologies may be generated, yielding conflicting phylogenetic information. It remains an open question how to best construct a consistent evolutionary history from conflicting phylogenetic signals<sup>6</sup>

Historically, reticulate processes were believed to occur at such a low frequency that they could be safely ignored when considering evolutionary relationships. However, new genomic data has shown that, particularly in microorganisms such as bacteria and archaea, reticulate processes are much more prevalent than originally expected [123]. Incompatibilities in the tree paradigm now appear as the rule, not the exception, which has led to calls for new representations of evolutionary relationships [45, 46]. Many have argued that, in light of new genomic evidence, the very notion of a universal tree of life must be discarded [94, 95]. This point has been argued most strongly by Ford Doolittle of Dalhousie University. In Figure 1.3, we see Doolittle's simplified representation of Tree of Life as it stands today, including only two of the most well-known reticulate events: the acquisition of mitochondria and chloroplasts from bacterial ancestors. We also see his representation of the Tree of Life as it would stand if additional large-scale reticulations were reflected – no longer is it clear that the tree is an appropriate metaphor.

Finally, reticulate evolutionary processes are of more than just historical interest for evolutionary studies, but play a substantial role in human health and disease. In HIV, frequent homologous recombination confounds our understanding of the epidemic's early and present history [25]. In influenza, segmental gene reassortments lead to antigenic novelty and the emergence of epidemics [118]. In several pathogenic bacteria, including *E. coli* and *S. aureus*, horizontal gene transfer has been responsible for the spread of antibiotic resistance genes [2, 42]. For example, the 2011 German *E. coli* outbreak was caused by a strain of *E. coli* that had acquired the ability to produce Shiga toxin [130].

---

<sup>6</sup>There exists a cottage industry of methods for aggregating conflicting *gene trees* into a consensus *species tree*, see [110].

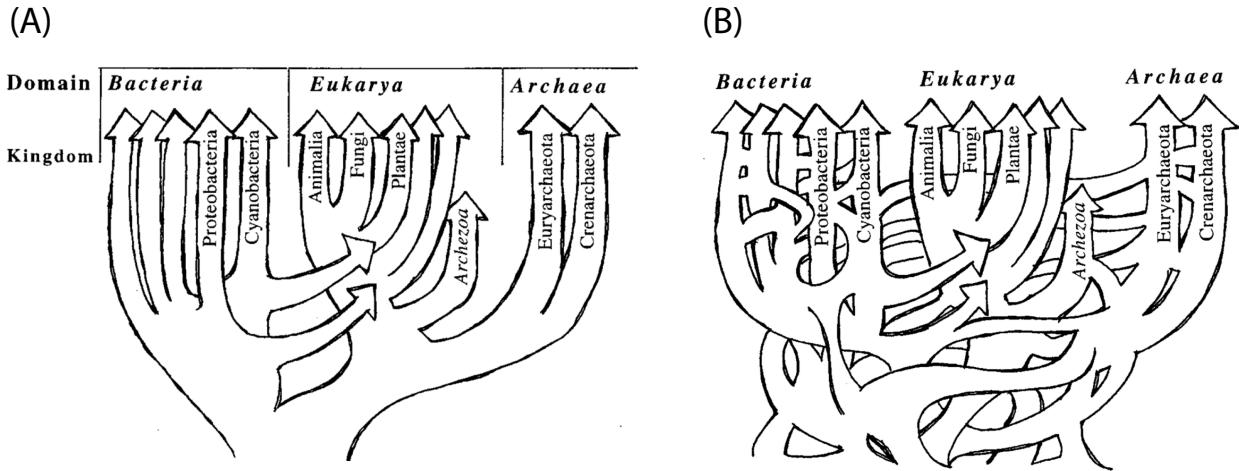


Figure 1.3: (A) W Ford Doolittle’s representation of the consensus universal tree of life. Only the most well known reticulations are reflected: the endosymbiosis of mitochondria and chloroplasts. (B) Doolittle’s speculative representation of the universal tree of life after accounting for reticulate evolution. While the three domains of life are still recognizable, patterns of divergence no longer follow a strictly treelike model. (From *Science*, vol. 284, issue 5423, page 2127. Reprinted with permission from AAAS.)

## 1.3 Evolution as a Topological Space

We propose the use of new computational techniques, borrowed from the field of applied topology, to capture and represent complex patterns of reticulate evolution.

Topology as a mathematical field is concerned with properties of spaces that are invariant under continuous deformation. Such properties can include, for example, connectedness and the presence of holes. Two objects are considered topologically equivalent if they can be deformed into one another without introducing any cuts or tears. As a paradigmatic example, consider the coffee mug and the donut (Figure 1.4). While seemingly different, it is not difficult to see that both objects consist of a single connected component that is wrapped around a single hole. Were the objects smoothly pliable they could be freely deformed into one another. Topologically, the two objects are equivalent.<sup>7</sup>

Algebraic topology quantifies our intuitive notions of shape using algebraic structures

---

<sup>7</sup>The two objects are topologically equivalent to a solid torus, which is represented as  $D^2 \times S^1$ , a solid two-dimensional disk wrapping around a circle.

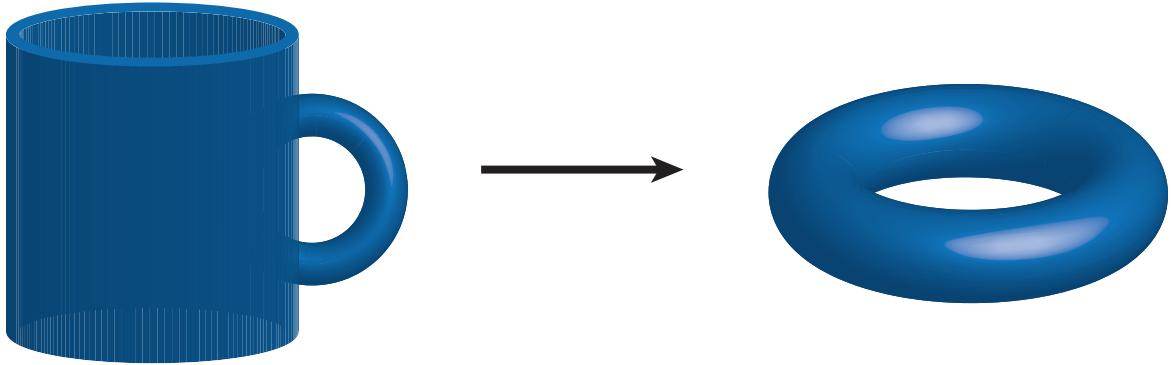


Figure 1.4: The paradigmatic example of topological equivalence. The coffee mug can be continuously deformed into the donut and are therefore topologically equivalent. Both exhibit the topology of a solid torus ( $D^2 \times S^1$ ).

to represent different invariants of a space. For our purposes, the most relevant invariants will be the *Betti numbers*. We give a more complete characterization of Betti numbers in Chapter 2, but the intuition is as follows. The Betti numbers are a collection of integers indexed by an integer  $n$  describing the connectivity of a space at different dimensions. First, we can think of  $\beta_0$  as representing the number of connected components, or clusters, in our space. Next, we can think of  $\beta_1$  as representing the number of one-dimensional loops in our space. Equivalently, this is the number of cuts needed to transform the space into something simply connected.<sup>8</sup> Higher Betti numbers,  $\beta_n$  for  $n > 1$  will correspond to higher dimensional holes. In our coffee mug example, because both objects have the same Betti numbers ( $\beta_0 = 1$ ,  $\beta_1 = 1$ , and  $\beta_n = 0$  for  $n > 1$ ), they are considered topologically equivalent. Our goal in this work will be to adopt a similar perspective and characterize evolutionary spaces as topological spaces using their Betti numbers.

To give a simple example, consider Figure 1.5. The example presents two possible scenarios describing the evolutionary relationships of three species, labeled  $a$ ,  $b$ , and  $c$ . For each

---

<sup>8</sup>In a simply connected space, any path between two points can be deformed into any other such path.

scenario, moving vertically up the object corresponds to moving backwards in time. Branch lengths correspond to evolutionary divergence. Internal vertices represent extinct ancestors of the three species, up to the root of the tree,  $r$ , which represents the most recent common ancestor. On the left, we have a simple tree topology relating the three species. Considering the shape of the tree, there is a single connected component, giving  $\beta_0 = 1$ . Further, we see that there are no loops formed by the branches, giving  $\beta_1 = 0$ . The object is therefore considered trivially contractible, a property which will hold for all tree topologies. On the right, we have a reticulate topology relating the three species. We can envision species  $b$  as being the reticulate offspring of parents ancestral to species  $a$  and  $c$ . That is, species  $b$  carries unique genetic material from both species  $a$  and species  $c$ . To account for this, two branches merge into the vertex that is directly ancestral to  $b$ . Considering the shape, there is again a single connected component, giving  $\beta_0 = 1$ . However, because of the reticulate event mixing material from  $a$  and  $c$ , there is now a loop formed in the topology, giving  $\beta_1 = 1$ . The object is no longer treelike and is characterized by a more complex topology. The Betti numbers capture the essential difference in the two evolutionary histories. Finally, we note that this is a conceptual example of how reticulate processes can be captured using topology – in practice, we do not have access to the true history, but must infer it from a finite sample.

Consider again Darwin’s branching phylogeny (Figure 1.1) and Doolittle’s modified representation after accounting for reticulate evolution (Figure 1.3). The two objects can be imagined to be representations of two different topological spaces. Darwin’s branching phylogeny is a tree and hence trivially contractible ( $\beta_n = 0$  for  $n > 0$ ). In contrast, Doolittle’s construction has a much more complex topology, with loops being formed where reticulate events have occurred. The object will be characterized by nonvanishing Betti numbers, the magnitude of which will be associated with the amount of reticulation that has occurred. The remainder of this thesis focuses on expanding this idea and applying it to real data sets with the goal of measuring the prevalence and scale of reticulate evolutionary events. Our aim will be to characterize reticulate exchange of genetic material by the parental sequences

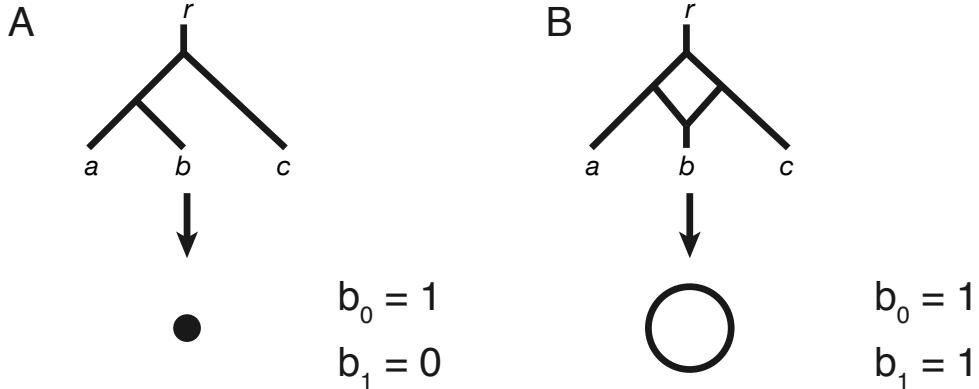


Figure 1.5: (A) A simple treelike phylogeny is contractible to a point. (B) A reticulate phylogeny that is equivalent to a circle and not contractible without a cut. The two spaces are not topologically equivalent and can be distinguished by their Betti numbers.

involved in the exchange, by the amount and identity of material exchanged (i.e., the genes or loci involved), and the frequency with which similar exchanges occur. Several important questions will be dealt with, such as how to construct topological spaces from finitely sampled sequence data, how to make comparisons among gene sets, and how to make statistical statements about reticulate events. We will address these questions by developing new techniques to construct and extract topological and statistical information from evolutionary data. In doing so, we provide a fuller understanding of evolutionary relationships than possible with current phylogenetic methods.

## 1.4 Thesis Organization

The remainder of this thesis is organized as follows.

In Chapter 2 we present background material on the topics discussed in this thesis. This discussion is chiefly structured into two pieces: (1) background on phylogenetics and population genetics, and (2) background on the methods we use from TDA.

In Part I, we develop two complementary approaches for analyzing sequence data using TDA. In Chapter 3, we propose methods of constructing topological spaces that generalize

standard constructions but are suited to the particular requirements of phylogenetic applications. We draw on previous work in phylogenetic networks and use homology to provide a quantitative assessment of reticulate processes. This work was published in [58]. In Chapter 4, we develop methods for performing statistical inference using summary statistics computed using methods from TDA. This is the first such use of TDA as a tool for performing parametric inference and should generalize to a wide range of application settings. This work was published in [59]

In Part II, we apply our approach to several problems in evolution and genomics. In Chapter 5 we study phages, viruses of single-celled microorganisms. We show how persistent homology recovers inconsistencies in existing morphology-based taxonomies, use a network approach to construct an alternative genome-based representation of phage relationships, and identify representative gene families conserved within phage populations. In Chapter 6 we study influenza, a common human pathogen. We show how persistent homology captures widespread patterns of reassortment, including nonrandom cosegregation of segments and barriers to subtype mixing. In contrast to traditional influenza studies, which focus on the phylogenetic branching patterns of only the two surface-marker proteins, we use Mapper combined with whole-genome data to represent influenza molecular relationships. We show unexpected relationships between divergent influenza subtypes. This work draws from results in [31] and [59]. In Chapter 7 we study pathogenic bacteria. We use two sources of data to measure rates of reticulation in both the core genome and the mobile genome across a range of species. Mapper is used to represent the population of *S. aureus* and analyze the spread of antibiotic resistance genes. The potential for the spreading of antibiotic resistance in the human microbiome is investigated. This work was published in [57]

Finally, in Chapter 8 we summarize these results and present future research directions.

# Chapter 2

## Background

This thesis uses newly developed approaches from applied topology to study problems in evolutionary biology and genomics. In this chapter we provide background material to motivate our approach. In Section 2.1 we introduce models of evolution and the types of genomic data we will consider. In Section 2.2 we provide a self-contained introduction to the primary methods of topological data analysis, including persistent homology and Mapper. Finally, in Section 2.3 we give simple examples of how the tools from TDA can be informative about reticulate evolution.

### 2.1 Evolutionary Biology and Genomics

In this section we present a basic introduction to molecular sequence data: what the data looks like, the processes by which it is generated, and the methods by which it is analyzed. Particular attention is paid to modes of reticulate evolution. Exposition for specific biological applications can be found in their respective chapters.

### 2.1.1 Genes and Genomes

The information required to express an organism's biological form and function is contained in the genome. At least one copy of the genome is packaged inside each cell of an organism. Physically, the genome is manifest as a polymer chain of nucleic acid, built on an alphabet of four nucleotide monomers: adenine, cytosine, guanine, and thymine. Abstractly, the genome is represented as a linear sequence of characters defined over the alphabet  $\{A, C, G, T/U\}$ .<sup>1,2</sup> Contained in this sequence are subsequences representing genes, which code for the protein products that ultimately affect function. Further embedded in the genome is a complex regulatory pattern of transcription factors controlling the expression of particular genes and directing cellular differentiation and development.

Following the central dogma of biology, DNA is transcribed into RNA, RNA is translated into amino acids, and amino acids are folded into proteins [39]. Proteins comprise the functional unit of biology.

Beyond simply coding for function, the genome includes an imprint of the evolutionary history that gave rise to the organism. By comparing the genomes of multiple organisms, inferences can be drawn about the evolutionary relationships among extant organisms as well as the processes that generated observed diversity. The field concerned with exploring these relationships is *comparative genomics*.

### 2.1.2 Evolutionary Processes

Evolution describes the gradual change in phenotypes arising from random variation and subject to natural selection. The processes giving rise to diversity can be classified into two types: clonal and reticulate.

---

<sup>1</sup>T in the case of DNA genomes. U in the case of RNA genomes.

<sup>2</sup>The linear representation can be misleading, as many organisms, primarily viruses and bacteria, have circular genomes.

### **2.1.2.1 Clonal Evolution**

Clonal evolution, or vertical evolution, is a process of self-reproduction whereby genetic material is transferred directly from parent to offspring. Population diversity is generated by stochastic mutation and maintained over multiple generations by random drift.

It is clonal evolution that Darwin had in mind when he described the idea of descent with modification, whereby a parent passes genomic information to an offspring subject to random drift. Importantly, because there is always a direct parent–offspring relationship, clonal evolution can be modeled with a binary tree model.

### **2.1.2.2 Reticulate Evolution**

Reticulate evolution, or horizontal evolution, refers to exchange or acquisition of genetic material via processes that do not reflect a direct parent–offspring relationship. As we will see, these processes can make inferences about historical evolutionary relationships difficult. Different types of reticulate processes occur in different types of organisms (summarized in Table 2.1).

Viruses replicate by infecting a host cell and then using the host cell machinery and resources to produce multiple copies of viral genetic material. The genetic material is then packaged into new virus particles which are shed off in order to infect new cells. Reticulation can occur when two virus particles coinfect the same host cell. During the replication process, genetic material can be exchanged in one of two ways: *reassortment* or *recombination* (the two processes are contrasted in Figure 2.1). Reassortment occurs in viruses whose genomes are segmented, such as influenza. Segments are similar to chromosomes, such that a single virus particle will contain a single copy of each segment. Coinfection of a single cell with two independent viruses results in packaging of segments taken from different virus particles. The result viral progeny will then be a genetic mixture of segments from each parental strain. Recombination, more common in non-segmented viruses such as HIV, involves a break-rejoin mechanism during the replication process. Here, an error in the polymerase during

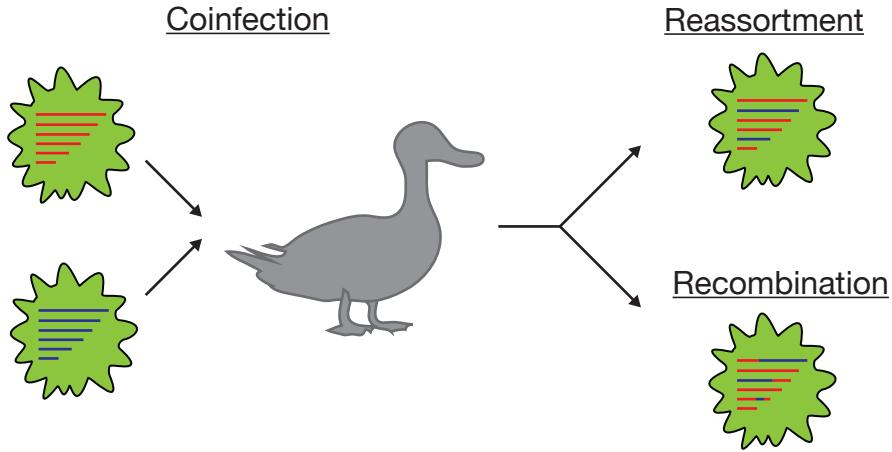


Figure 2.1: The two modes of viral reticulation. Coinfection of the same host cell can lead to either reassortment, in which whole viral segments are exchanged, or recombination, in which breakpoints can occur within segments. The former process is common in influenza, the latter in HIV. The end result, however, is a novel virus particle which shares genetic information from both parents.

replication can result in an incomplete copy of the genome (a break). At this point, several cellular processes involved in repair can be recruited to complete the replication process using a homologous region. If coinfection has occurred, it is possible for these processes to initiate repair using material from a different parental strain. The outcome will be novel genetic material that includes a crossover from one strain to another. Break-rejoin crossover is a type of *homologous recombination*.

In bacteria and other prokaryotes, reticulate evolution can occur when foreign DNA from a donor is acquired by a target organism and integrated into its genome. Three generic mechanisms have been identified, depending on the route by which foreign DNA is acquired [123]:

1. *Conjugation*. Direct cell-to-cell contact between donor and recipient resulting in transfer of plasmid.
2. *Transformation*. Foreign DNA acquired via uptake from freely circulating DNA in the environment.

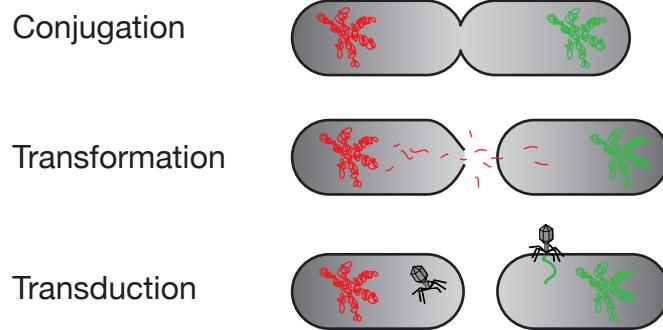


Figure 2.2: Three modes of viral reticulation. (1) Conjugation, in which direct cell-to-cell contact results in transfer of genetic material; (2) Transformation, in which foreign DNA is acquired via uptake from freely circulating DNA in the environment; and (3) Transduction, in which exchange of genetic material is mediated by a virus or phage particle.

### 3. *Transduction*. Virus-mediated transfer for foreign DNA from an infected donor cell.

A visualization of these three mechanisms is shown in Figure 2.2. Because these mechanisms can often lead to the acquisition of novel sequences coding for genes not in the recipient organism, reticulate evolution in prokaryotes is often called *horizontal gene transfer* or *lateral gene transfer*.

In eukaryotes, several reticulate processes have been identified. We mention two such processes: hybrid speciation and meiotic recombination. These two processes act at very different scales, however the outcome is the same: a unique offspring with genetic material drawing from both parents.

First, hybrid speciation refers to the cross-breeding of animals or plants of different species. This mixing of genetic material can lead to the development of offspring with a phenotype distinct from both parents. Hybrid speciation was originally believed to be a rare occurrence in nature and hybrid offspring to be infertile. However, recent genomic data has demonstrated that hybridization occurs quite frequently in plants [4, 5].

Second, meiotic recombination refers to a specialized process for generating diversity that occurs in sexually-reproducing polyploid organisms, such as humans, during meiosis. Meiosis is the process by which a single cell containing  $n$  copies of each chromosome results

Table 2.1: Reticulate processes in biology across kingdoms

Organism	Process	Description
Virus	Reassortment	Exchange of discrete genomic segments
	Recombination	Intragenomic homologous crossover
Bacteria	Transformation	Acquisition of foreign DNA in environment
	Transduction	Viral-mediated exchange
	Conjugation	Cell-to-cell contact and exchange
Eukaryotes	Meiotic Recombination	Homologous crossover during meiosis
	Hybrid Speciation	Fertilization across species boundaries

in four distinct cells each with  $n/2$  copies of each chromosome. These special cells are called gametes. Sexual reproduction consists of the fusion of two gametes during fertilization to form a zygote, which ultimately develops into a viable offspring. Meiosis is a multi-step process consisting of an initial round of DNA replication followed by two rounds of cell division. Meiotic recombination occurs after the initial round of DNA replication and prior to cell division. After DNA replication, there are two copies of each homologous chromosome that are joined at a centromere. The two sets of chromosomes then pair with each other and exchange DNA through physical interactions known as crossovers.<sup>3</sup> This is another example of homologous recombination and results in new allelic patterns mixing genetic information from both parents.<sup>4</sup> After crossover occurs, two phases of cellular division result in gametes with  $n/2$  copies of each chromosome.

The presence of reticulate processes in a set of organisms can be most clearly identified by comparing phylogenetic relationships built from different genomic segments. A general practice is to construct the set of *gene trees* which reflect ancestral branching patterns at specific loci. If a reticulate event has occurred, it implies that the branching patterns of different genes will not agree. A subfield of comparative genomics is concerned with building *species trees* from sets of gene trees [110].

---

<sup>3</sup>These crossovers have been shown to occur nonrandomly at recombination hotspots regulated by binding motifs for the PRDM9 protein [10, 26].

<sup>4</sup>Patterns of shared alleles define the concept of *linkage*.

However, in the case where there is substantial disagreement among gene trees, the very notion of a species tree may be flawed. Traditionally, evolutionary biology has concerned itself with characterizing relationships in light of vertical evolution alone. However, increasing evidence has pointed to the important role played by horizontal evolution, particularly in prokaryotic evolution [71, 70]. Between 10% to 16% of the *E. coli* genome is believed to have arisen from horizontal gene transfer [123].

### 2.1.3 Mathematical Models of Evolution

Mathematical population genetics is concerned with properties of populations as they are subject to evolutionary forces over long time scales. These forces include natural selection, genetic drift, mutation, and recombination. Historically the input data for population genetics models was comparative studies of allele frequencies across populations. These studies have primarily been replaced by large-scale genomic surveys which have provided unprecedented insight into ancient population structure and historical migrations.

These models allow scientists to two things: (1) simulate genomic data under realistic processes and (2) build statistical models to estimate biological parameters from data.

#### 2.1.3.1 The Wright-Fisher Model

The Wright-Fisher model is a forward time simulation of an evolving population. In the simplest case, the model describes neutral evolution of a constant population size with no structure and constant genome length. The model proceeds in units of generations. At each generation, a member of the population is an offspring of a randomly selected ancestor from the previous generation. This offspring inherits its ancestors genomes, with mutations introduced at some base rate  $\mu$ . A member of previous generation with no offspring will be considered extinct.

### 2.1.3.2 The Coalescent Process

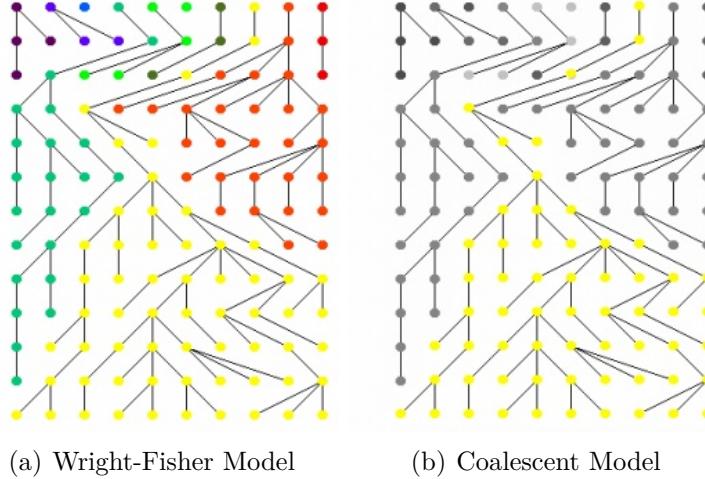
The coalescent process is a stochastic model that generates the genealogy of individuals sampled from an evolving population [148]. The genealogy is then used to simulate the genetic sequences of the sample. This model is essential to many methods commonly used in population genetics. Starting with a present-day sample of  $n$  individuals, each individual's lineage is traced backward in time, towards a mutual common ancestor. Two separate lineages collapse via a coalescence event, representing the sharing of an ancestor by the two lineages. The stochastic process ends when all lineages of all sampled individuals collapse into a single common ancestor. In this process, if the total (diploid) population size  $N$  is sufficiently large, then the expected time before a coalescence event, in units of  $2N$  generations, is approximately exponentially distributed:

$$P(T_k = t) \approx \binom{k}{2} e^{-(\frac{k}{2})t}, \quad (2.1)$$

where  $T_k$  is the time that it takes for  $k$  individual lineages to collapse into  $k - 1$  lineages.

After generating a genealogy, the genetic sequences of the sample can be simulated by placing mutations on the individual branches of the lineage. The number of mutations on each branch is Poisson-distributed with mean  $\theta t/2$ , where  $t$  is the branch length and  $\theta$  is the population-scaled mutation rate. In this model, the average *genetic distance* between any two sampled individuals, defined by the number of mutations separating them, is  $\theta$ .

The coalescent with recombination is an extension of this model that allows different genetic loci to have different genealogies. Looking backward in time, recombination is modeled as a splitting event, occurring at a rate determined by population-scaled recombination rate  $\rho$ , such that an individual has a different ancestor at different loci. Evolutionary histories are no longer represented by a tree, but rather by an *ancestral recombination graph*. Recombination is the component of the model generating nontrivial topology by introducing deviations from a contractible tree structure, and is the component which we would like to quantify. Coalescent simulations were performed using `ms` [81].



(a) Wright-Fisher Model      (b) Coalescent Model

Figure 2.3: Two models for simulating evolutionary data. On the left, the Wright-Fisher model simulates a sample of  $n$  individuals in the forward direction. At each generation,  $n$  offspring choose a parent from the previous generation at random. After  $t$  generations, some initial lineages will have died off, while others will become dominant in the population. On the right, the coalescent model simulates the sample in the reverse direction. At each reverse generation, individuals merge, or coalesce, with some probability, until they reach a single most recent common ancestor (MRCA). The intuition behind the approach is that lineages that have gone extinct will not contribute to the present day observed diversity, are therefore inaccessible, and do not need to be simulated. This approach reduces the data that needs to be simulated and increases the computational performance of the models.

### 2.1.3.3 Metrics on Sequences

Evolutionary models require a notion of genetic divergence between sequences. This leads to a discussion of the types of metrics that can be put on sets of sequences.<sup>5</sup>

The simplest model, and the one most commonly adopted in this thesis, is the Hamming metric, which simply counts the proportion of sites that differ between two aligned sequences. For example, for two sequence  $s_1 = ACTTGAC$  and  $s_2 = AAGTGGC$ ,  $d_H(s_1, s_2) = 3/7$ . In general, the Hamming metric will underestimate divergences by not accounting for the

---

<sup>5</sup>Before sequences can be compared, they must first be *aligned*. A sequence alignment arranges the characters in a set of sequences into columns such that individual characters sharing an evolutionary identity are in the same column. Alignment is necessary because random insertion and deletion of nucleotides can change the relative positions of related bases. The difficulty of performing an alignment will largely depend on the amount of evolutionary divergence in the set of sequences under consideration. Sequence alignment is a well studied topic but largely beyond the scope of this thesis, where we assume sufficient sequence similarity such that alignment can be performed with high confidence.

possibility of back mutations.<sup>6</sup>

More biologically motivated models will introduce corrections to account for assumptions about how sequences evolve. These assumptions include the base frequency of each nucleotide as well as the substitution rates for each type of mutation. The simplest of these models is the *Jukes-Cantor model*. This model defines an equal substitution rate  $\mu$ . Inverting the probability of an alteration gives the divergence. The Jukes-Cantor metric is defined as

$$d_{JC} = -\frac{3}{4} \ln(1 - \frac{4}{3}p), \quad (2.2)$$

where  $p$  is the proportion of sites that are different.

#### 2.1.4 Phylogenetic Methods

A phylogenetic tree is a binary tree in which leaves are associated with particular species or taxa, and the branching pattern of the tree reflects diverging evolutionary relationships. Branch lengths on the tree are associated with evolutionary divergence between sets of taxa. A tree can be either rooted, in which case a particular point on the tree is identified as the most recent common ancestor and the temporal order of branching is inferred, or unrooted, in which case only the branching pattern is represented but no statements about their temporal order are inferred. Typically sequence data alone is not sufficient to root a tree – an estimate of the mutation rate under an evolutionary model is also required. See Figure 2.4 for an example of the two types of trees. In this work we primarily deal with unrooted trees.

Molecular phylogenetics refers to a large collection of methods for inferring branching patterns from aligned molecular sequence data.<sup>7</sup> In general, the problem of finding an optimal tree associated with sequence data is NP-complete [66], however several approximate methods have been developed. The primary types of methods include maximum parsimony, distance-matrix methods, maximum likelihood (ML), and Bayesian inference. Maximum

---

<sup>6</sup>A double mutation of the form  $A \rightarrow C \rightarrow A$ .

<sup>7</sup>See Felsenstein's *Inferring Phylogenies* for a readable and thorough introduction to the field [64].

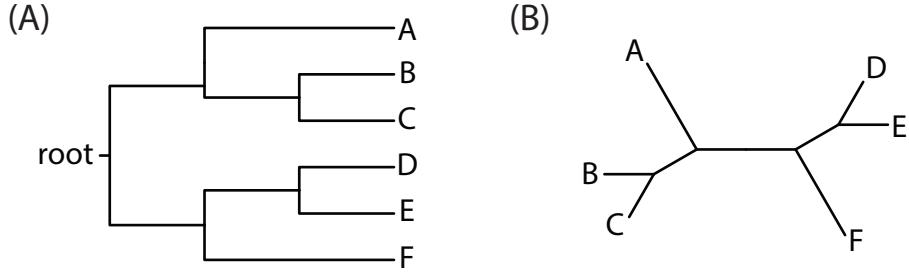


Figure 2.4: (A) A rooted tree and (B) an unrooted tree on six leaves. In a rooted tree, a particular point on the tree is identified as the most recent ancestor. Time is measured along the horizontal axis. In an unrooted tree, only the pattern of divergence is represented. From sequence data, often times only an unrooted tree can be inferred.

parsimony attempts to find the phylogenetic tree that minimizes the number of evolutionary changes required to explain the observed sequences. Distance-matrix methods first compute a matrix of pairwise distances between taxa and then find the tree that best approximates these distances. ML and Bayesian methods use specific models of evolution to assign probability distributions over trees. In this work we concentrate on distance-matrix methods because of their close connection with the finite metric spaces considered in applied topology.

#### 2.1.4.1 Distance-Matrix Methods

Given a set of aligned molecular sequences, distance-matrix methods first compute the pairwise matrix of genetic distances using one of the metrics as described in Section 2.1.3.3. Then, the binary tree that best approximates those distances is iteratively fit to this data. This approach to phylogenetic inference were introduced by Cavalli-Sforza and Edwards in 1967 [30] and Fitch and Margoliash in 1967 [65]. The Fitch-Margoliash method uses a weighted least squares approach to tree-fitting, such that larger distances are weighted less, due to higher chances for random error. Distance-matrix methods are popular for their high speed and scalability as well as high accuracy in most cases.

**Data:**  $n \times n$  distance matrix  $D$

**Result:** Phylogenetic tree on  $n$  leaves

**while** Tree not fully resolved ( $n > 3$ ) **do**

    Compute  $Q$  matrix:

$$Q(i, j) = (n - 2)d(i, j) - \sum_{k=1}^n d(i, k) - \sum_{k=1}^n d(j, k);$$

    Identify pair of taxa  $i, j$  that minimizes  $Q(i, j)$ ;

    Create new interior node  $u$  that joins  $i$  and  $j$  with edge length:

$$D(i, u) = \frac{1}{2}D(i, j) + \frac{1}{2(n-2)} [\sum_{k=1}^n D(i, k) - \sum_{k=1}^n D(j, k)];$$

$$D(j, u) = D(i, j) - D(i, u);$$

    Create new  $(n - 1) \times (n - 1)$  distance matrix where:

$$D(u, k) = \frac{1}{2}[d(f, k) + d(g, k) - d(f, g)];$$

**end**

**Algorithm 1:** The Neighbor Joining Algorithm. Adapted from [153]

Currently, the most widely implemented distance-matrix method is neighbor-joining.<sup>8</sup>

One particular reason neighbor-joining is popular is that under certain conditions, discussed below, it has been shown to exactly recover the correct tree. The neighbor-joining algorithm is a greedy approach to tree construction that iteratively joins the two closest nodes until a tree is fully resolved. The neighbor-joining algorithm is described in Algorithm 1.

#### 2.1.4.2 Additive Metrics and the Four Point Condition

Arbitrary distance matrices are unlikely to admit a tree representation. Those that do are called *additive metrics*, because they can be represented as an additive tree. Additivity is the property that the distance between any two nodes will be equal to the sum of the branch lengths between them. A distance matrix admits a tree representation if and only if it is additive.

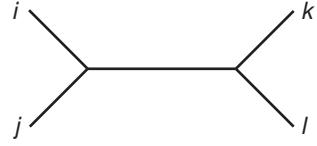
There is a straight-forward condition that must be satisfied for additivity, known as the *four point condition*. For a distance matrix to admit a tree representation,

$$d_{ij} + d_{kl} \leq \max\{d_{ik} + d_{jl}, d_{il} + d_{jk}\} \quad (2.3)$$

for any four nodes  $\{i, j, k, l\}$ . The condition implies that there is a labeling on the four nodes

---

<sup>8</sup>Neighbor joining was introduced by Saitou and Nei in 1987 [133].



$$d_{ij} + d_{kl} \leq d_{ik} + d_{jl} = d_{il} + d_{jk}$$

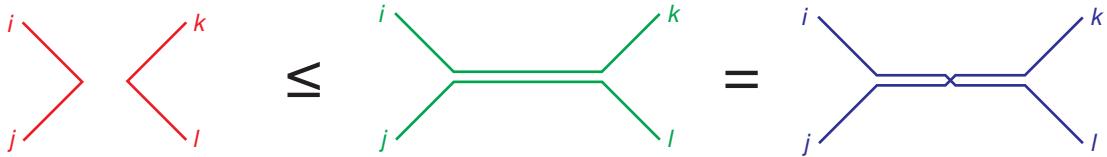


Figure 2.5: A visual interpretation of the four point condition for additivity. For any four leaves, there exists a labeling  $\{i, j, k, l\}$  such that  $d_{ij} + d_{kl} \leq d_{ik} + d_{jl} = d_{il} + d_{jk}$ . Of the three possible ways of arranging the sums of distances, two will involve traversing the internal branch, while one will involve only external branches.

such that

$$d_{ij} + d_{kl} \leq d_{ik} + d_{jl} = d_{il} + d_{jk}. \quad (2.4)$$

A visual interpretation of this condition is shown in Figure 2.5.

Sequence data can fail to be additive for several reasons. First, sequencing error. Errors can introduce noise into the measured genetic distances. Second, homoplasy. A homoplasy occurs when the same mutation is introduced multiple times in a set of organisms. The presence of homoplasy will underestimate genetic distance between taxa. Third, reticulate evolution. As described previously, in cases of reticulate evolution no tree will accurately describe the observed data. In this case, one can either attempt to find the tree that best fits the data, or search for an alternative representation of phylogenetic relationships.

#### 2.1.4.3 Number of Tree Topologies

Labeled trees on a fixed set on set of leaves can be distinguished by their topology, which refers to the arrangement of leaf labels corresponding to a particular evolutionary history.<sup>9</sup> The number of unrooted bifurcating tree topologies with  $L$  leaves is  $\mathcal{T}(L) = (2L - 5)!!$ .<sup>10</sup> This can be easily shown using induction. For  $L = 3$ , we have  $\mathcal{T}(3) = 1$  and 3 branches. To pass to  $L = 4$ , we can add the fourth leaf to any of the 3 branches, resulting in 3 different topologies. For  $L = 4$ , we have  $\mathcal{T}(4) = 3$ . Every time we add a leaf, we add two branches – one external and one internal. For  $L = n$ , we have  $\mathcal{T}(n) = (2n - 5)!!$  and  $2n - 3$  branches. For  $L = n + 1$ , we can add the new external branch to any of the current  $2n - 3$  branches. A rooted tree with  $L$  leaves can be considered as an unrooted tree with  $L + 1$  leaves. Therefore, the number of rooted bifurcating tree topologies with  $L$  leaves is  $(2L - 3)!!$  As can be seen, the number of tree topologies explodes with the number of leaves.<sup>11</sup> See Figure 2.6.

#### 2.1.4.4 The Space of Phylogenetic Trees

An unrooted phylogenetic tree with  $L$  leaves is characterized by its topology and the lengths of each branch. As shown in the previous section, there are  $(2L - 5)!!$  possible unrooted topologies. There are  $2L - 3$  total branches, of which  $L$  are external branches and  $L - 3$  are internal branches. Tree spaces refers to an abstract construction for representing each possible tree as a point in a geometric space. These studies were initiated by Andreas Dress and colleagues, who introduced a formalism known as *T-theory* (see [52, 51, 49]). We give here a brief flavor of these ideas; additional exposition can be found in [125, §7].

Consider a set of  $L$  leaves. A dissimilarity map is defined on  $L$  as  $\delta : L \times L \rightarrow \mathbb{R}$ , where  $\delta(l, l) = 0$  and  $\delta(l, m) = \delta(m, l)$ . There are  $\binom{L}{2}$  distances; the set of dissimilarity

---

<sup>9</sup>The use of the term topology here is standard in phylogenetics, but distinct from that in mathematical topology.

<sup>10</sup>The double factorial is defined as  $n!! = n(n - 2)(n - 4) \dots$ .

<sup>11</sup>As was observed by Walter Fitch, for 22 species there are on the order of Avogadro's number of topologies. ( $N_{22} = 3.20e23$ ,  $N_A = 6.02 \times 10^{23}$ )

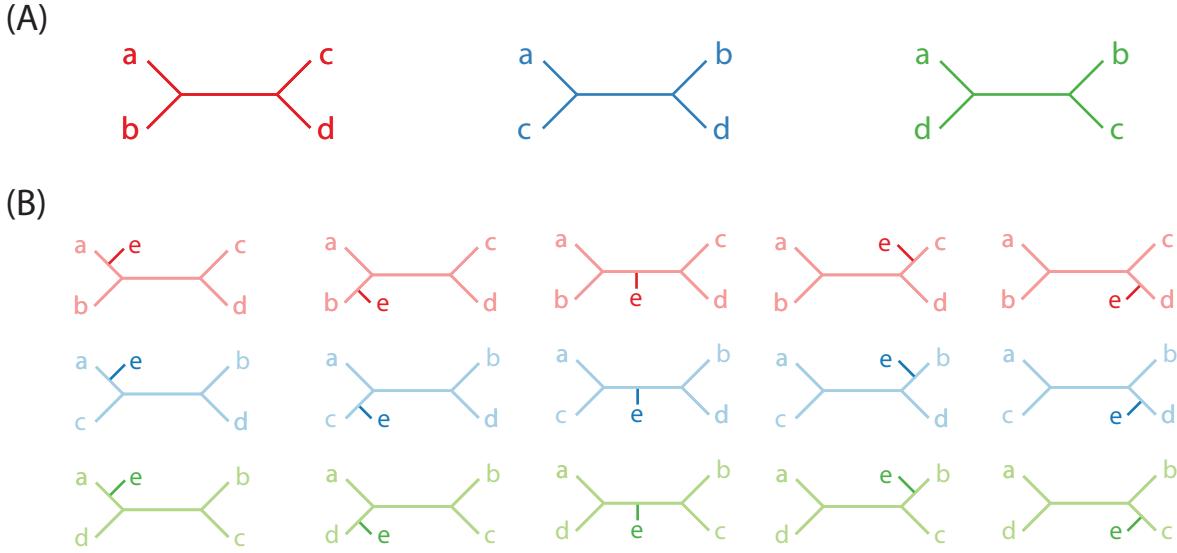


Figure 2.6: Enumerating tree topologies on labeled sets of leaves. (A) There are three unique tree topologies on four leaves. (B) There are fifteen distinct tree topologies on five leaves. Inductively, the fifth leaf can be added as a branch to each branch. In general, there are  $\mathcal{T}(L) = (2L - 5)!!$  topologies on  $L$  leaves.

maps forms a vector space of dimension  $\binom{L}{2}$ . Furthermore, the set of all metrics will be the subspace of  $\mathbb{R}^{\binom{L}{2}}$  that satisfies the triangle inequality. The space of trees is defined as the set  $\mathcal{T}_L$  of dissimilarity maps that satisfy the four-point condition. The space can be logically decomposed into subspaces corresponding to a particular choice of topology. This will be taken as the union of  $(2L - 5)!!$  subspaces, each of dimension  $2L - 3$ . Each subspace will have the structure of a metric cone in the space  $\mathbb{R}^{\binom{L}{2}}$ .

The geometric structure of this space was carefully studied by Billera, Holmes, and Vogtmann (BHV) in [15]. In that paper, the authors specifically considered rooted trees with zero-length external branches, a space denoted as  $\text{BHV}_L$ , but the basic intuition generalizes to other types of trees. They defined a geodesic distance between trees of different topology and used it to define various metric properties on tree space. This analysis was extended by Zairis *et al.* in [161], in which unrooted trees with non-zero external branches were considered. The external branches are constrained to sit in the positive open orthant  $(\mathbb{R}^{\geq 0})^L$ .

An evolutionary moduli space is then defined as the product

$$\Sigma_L = \text{BHV}_{L-1} \times (\mathbb{R}_{\geq 0})^L. \quad (2.5)$$

The tree space construction allows one to define statistics, such as means and variances, on collections of trees in a meaningful way.

We show an example of the tree space construction on  $L = 4$  and  $L = 5$  leaves in Figure 2.7. The case of  $L = 4$  is particularly simple to analyze. The metric cone is a subspace of  $\mathbb{R}^{\binom{4}{2}=6}$ . There are  $(2*4-5)!! = 3$  tree topologies, corresponding to the patterns  $((a, b), (c, d))$ ,  $((a, c), (b, d))$ , and  $((a, d), (b, c))$ . There are  $(2 * 4 - 3) = 5$  branches: each topology will be a subspace in  $\mathbb{R}^5$ . The intersection of the subspace of each topology is a space in  $\mathbb{R}^4$ . The case of  $L = 5$  also has a relatively simple structure. There are fifteen possible topologies, each with two internal branches. Each topology forms a hyperplane of dimension  $\mathbb{R}^7$ . Combinatorially, the topologies can be arranged as a Petersen graph. Intersections of three hyperplanes will correspond to degenerate cases with one internal branch is not resolved, as shown in Figure 2.7B. These facets sit in  $\mathbb{R}^6$ . It is important to think of the entire Petersen structure as being a cone, the origin of which is the 5-dimensional subspace consisting of only external leaves (see [15, Figure 14]).

Naturally, most data will not sit in  $\mathcal{T}$ . Whether or not this is simply due to noise or reflects reticulate processes will depend on the particular dataset. We can view the goal of phylogenetics as finding the best tree projection  $\delta_T \in \mathcal{T}$  for arbitrary metric data  $X$ .

#### 2.1.4.5 Phylogenetic Networks

There are several existing methods for representing reticulate evolution. Most of these methods generalize phylogenetic trees into *phylogenetic networks*, which attempt to reconcile the presence of horizontal evolution in sequence data. However, most simply present corrections to phylogenetic trees, which can fail in cases where horizontal evolution is pervasive, as in many prokaryote datasets. Additionally, the resulting networks can be complex and difficult to interpret quantitatively. This can make it difficult to distinguish between phylogenetic

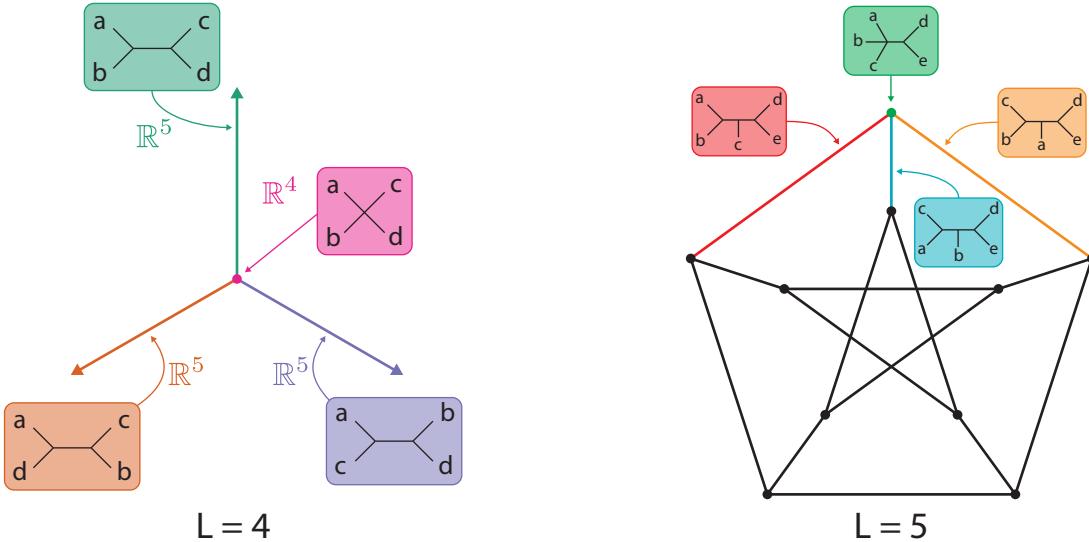


Figure 2.7: Examples of geometric representations of the space of trees on  $L = 4$  and  $L = 5$  leaves. (A) On four leaves the metric cone is a subspace of  $\mathbb{R}^6$ . There are three tree topologies, each of which corresponds to a 5-dimensional cone inside. The three topologies share a  $\mathbb{R}^1$  facet corresponding to the degenerate topology.(B) On five leaves the metric cone is a subspace of  $\mathbb{R}$ . There are fifteen tree topologies, each of which corresponds to a 7-dimensional cone. The geometric structure of the space will map to a *Petersen graph*, as shown. There are 10 degenerate cases in which one internal branch is not resolved; these correspond to 6-dimensional facets, each joining three distinct topologies. The  $n = 5$  subfigure is an adaptation of Figure 3.5 in [125, Ch 3].

incompatibilities due to noisy sampling and due to true reticulations. An example of a phylogenetic network using the split network approach is shown in Figure 2.8. Other methods include neighbor-net and median networks. Techniques such as phylogenetic networks and ancestral recombination graphs have been developed to describe reticulate evolution, but they have had only limited success due to difficulties of biological interpretation and computational infeasibility in all but the smallest datasets.

## 2.2 Topological Data Analysis

Topology is the branch of mathematics that formalizes our intuitive notions of shape. More concretely, topology provides the methods to characterize the properties of objects and

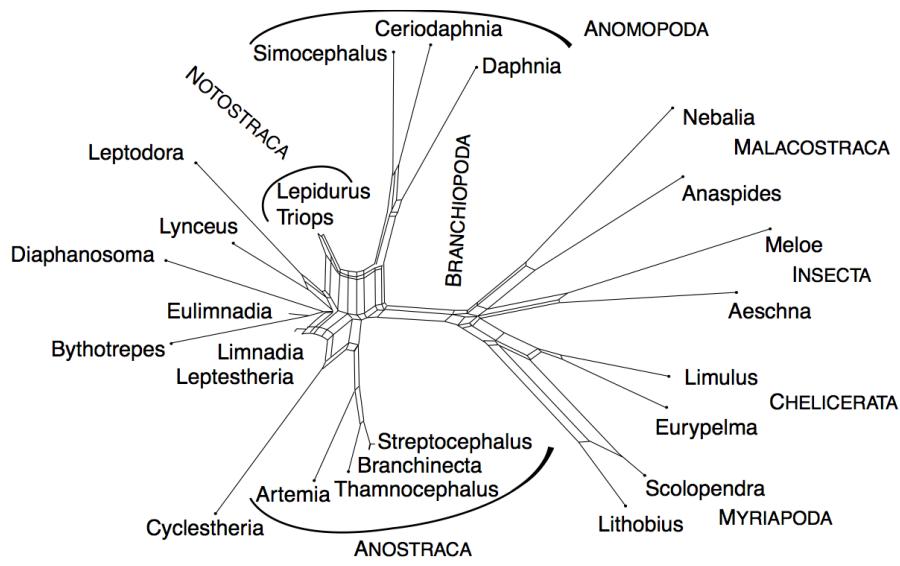


Figure 2.8: Example of a split network of genus Branchiopoda and outgroups. Computed using the Neighbor-Net algorithm. Phylogenetic incompatibilities are represented by conflicting splits. reprinted from BMC Evolutionary Biology 7:147 (2007).

spaces that remain invariant under continuous deformation. For example, transforming a circle into an ellipse by compressing along one axis does not change the fact that the object encloses a single loop. Or, as we saw in the introduction, the coffee mug can be continuously deformed into the donut. Likewise, if we take a tree and change the lengths of its branches, the tree remains a tree.<sup>12</sup> In each of these examples, while the deformation has substantially altered *local* properties of the space, on a *global* level certain essential characteristics have remained unchanged. From the perspective of topology, the spaces are considered identical. The question topology addresses is how to formalize the idea of global shape in order for it to be reasoned about systematically.

Algebraic topology solves this problem by associating to objects certain algebraic objects (an integer, for instance) that do not change under continuous deformation. These objects capture properties like the number of connected components, the number of loops, or the number of holes in an object, and represent *topological invariants* of a space. Two spaces can only be deformed into one other if they share the same invariants. For example, the circle and ellipse are identified as equivalent by the presence of a single loop. Neither can be deformed into a tree without introducing a cut, which would be a discontinuous deformation. Using these invariants, powerful ideas from abstract algebra can be used to manipulate and reason about shape.

While topology has traditionally developed through the study of abstract spaces, leading to very rich and beautiful constructions<sup>13</sup>, data does not come in the form of perfect continuous spaces. Recent effort over the past 15 years has focused on developing methods to apply topology to real world problems in science and engineering. This work, collectively falling under the heading of *topological data analysis* (TDA), has focused on efficient algorithms for computing topological invariants from finite, noisy data. TDA now encompasses a wide

<sup>12</sup>As was mentioned in Section 2.1.4.3, it is important to draw a distinction between the notion of tree topology, in which the branch patterns determines the topology, and global topology, in which all trees are equivalent. While the former is more common in the phylogenetics community, here we consider the latter.

<sup>13</sup>For example, see the work of Thurston on low-dimensional topology

range of efforts and can now be considered a branch of applied mathematics in its own right. It has emerged from substantial interdisciplinary effort between mathematicians, computer scientists, and domain experts.

In practice, a typical workflow for applying TDA to data is as follows. Data comes in the form of a set of  $n$  observations with  $p$  attributes, where  $p$  is often very large. The data is assumed to be a finite sample from a more complex space, from which we wish to infer either global structure or an underlying model. The data is represented as a finite point cloud: a set of  $n$  points in  $p$  dimensions with a notion of distance. The point cloud is transformed into a discrete topological space by associating different sets of points with each other, forming essentially a higher-dimensional analog of a graph. The associations can be constructed in different ways – for instance, one of the simplest constructions associates points within a certain distance  $d$  from one another. Computational approaches are then used to measure informative topological properties from the space.

In this thesis, we use methods from TDA to study problems in evolutionary biology and genomics. Our data is typically aligned genomic sequences from sets of related organisms, where features are the residues at each site. If our sequences are each of length  $L$ , then we can imagine our data as points in an  $L$ -dimensional sequence space. A genetic sequence metric, such as the Hamming metric, measures distance.

The two main methods from TDA that we employ are *persistent homology* and *Mapper*. Persistent homology provides a way to efficiently compute the topological invariants of a space across multiple scales, while Mapper provides an approach for condensed representation and visualization of high-dimensional data. In this section, we provide an overview and discussion of these two methods from the perspective of an end-user, treating each method as a pipeline for transforming from raw data to a concise topological summary. While the mathematical literature on these methods is extremely deep, our goal is to explain things in sufficient detail for a wide audience to grasp the main ideas. We therefore include a brief introduction of the basic mathematical concepts we employ. The primary concept we require

is *homology*, a particular way in which topological invariants can be assigned to spaces.

The following sections draw on several excellent reviews of TDA, including [27], [56], and [69]. A more thorough introduction to algebraic topology can be found in [75].

### 2.2.1 Preliminaries

As stated above, our data is a set of  $n$  points,  $S = \{s_1, \dots, s_n\}$ . Each point is a vector with  $p$  features,  $s_i = (s_{i1}, \dots, s_{ip})$ . We refer to the collection of points, embedded in a space with an appropriate metric structure, as a point cloud. We wish to associate a collection of algebraic objects to the point cloud in order to quantify its shape. To do so, our first step is to construct a topological structure on top of the point cloud, called a *simplicial complex*. The structure will consist of a set of simplices pieced together in such a way that they approximate the shape of the point cloud. Shape is then quantified using *homology*. This section provides the definitions necessary to understand homology.

#### 2.2.1.1 Simplices and Simplicial Complexes

The building blocks of our topological structures are simplices. A *simplex* is something like a point, a line, a triangle, or any higher-dimensional generalization of such. Formally, a  $k$ -simplex is a  $k$ -dimensional polytope which is the convex hull of  $k + 1$  vertices, as shown in Figure 2.9. A simplex can be represented by its list of vertices, i.e.  $\sigma = (s_1, s_2, s_3)$ . An  $m$ -face of a simplex is the space spanned by the set of  $m + 1$  vertices, and is itself a simplex. For example, the 0-faces and 1-faces of a simplex are its vertices and edges, respectively. The  $(k - 1)$ -faces (faces of co-dimension 1) of a  $k$ -simplex are called facets. Facets are represented as  $\sigma_{(-i)}$ , which implies the facet generated by elimination of the  $i$ -th vertex.

A *finite simplicial complex*  $K$  is built on the vertex set  $S$  from simplices glued together in such a way that (1) any face of a simplex in  $K$  is also in  $K$ , and (2) the non-empty intersection of any two simplices in  $K$  is a face of both simplices. An example of a simplicial complex is shown in Figure 2.10.

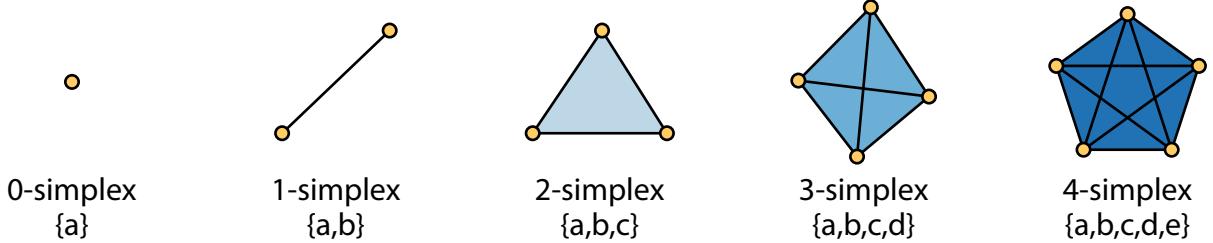


Figure 2.9: Simplices are the fundamental building blocks of our topological structures. They can be thought of as triangles generalized to arbitrary dimension. Here we show  $k$ -simplices for  $k = 0$  to  $k = 4$ .

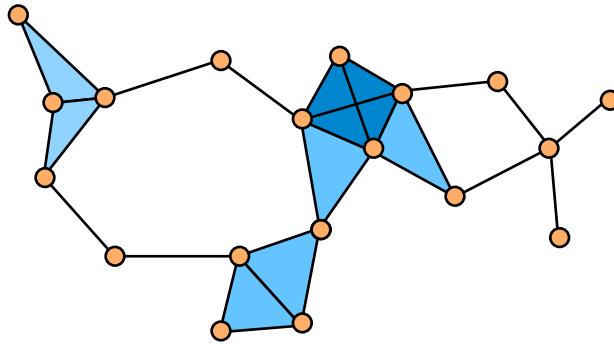


Figure 2.10: A finite simplicial complex  $K$  is an object built from a finite number of simplices, glued together in such a way that (1) any face of a simplex in  $K$  is also in  $K$ , and (2) the non-empty intersection of any two simplices is face of both simplices.

In order to compute homology and formally define the notion of holes, we need to define certain combinatorial operations that can be performed on a simplicial complex  $K$ . In general, these operations will act on subsets of simplices of fixed dimension  $k$ . These subsets are called  *$k$ -chains*, and can be represented as formal sums  $C_k = \sum_j \alpha_j \sigma_j$ . The coefficients  $\alpha_j$  will be taken over  $\mathbb{Z}_2$  (i.e. 0 and 1). Two consequences of this choice are (1)  $\sigma + \sigma = 0$ , and (2) we consider simplices without regard to orientation.<sup>14</sup>

An important operator is the boundary operator,  $\partial : C_k \rightarrow C_{k-1}$ . The boundary of a

---

<sup>14</sup>In general, an algebraic topology can be defined with coefficients in arbitrary fields. We use  $\mathbb{Z}_2$  for simplicity, efficiency, and because properties, such as torsion, that arise over more complex fields are not expected to be present in the biological data we consider. It is important to keep this in mind, as it was in fact shown that torsion can arise in real data in [29]. In that paper, an association was shown between the space of natural images and the Klein bottle.

simplex  $\sigma$ ,  $\partial_k \sigma$ , is the sum of its facets.

$$\partial_k \sigma = \sum_i \sigma_{(-i)} \quad (2.6)$$

The boundary of a chain is  $\partial C = \sum_j \partial \sigma_j$ . As a simple example, consider the 2-simplex  $\Delta$  defined by vertices  $\Delta = (a, b, c)$ . We have  $\partial \Delta = (a, b) + (b, c) + (a, c)$ . Further, we have  $\partial \partial \Delta = 2(a) + 2(b) + 2(c) = 0$ . In fact, the property  $\partial \partial C = 0$  will hold for any chain  $C$ .

We can additionally define more refined chains. A *cycle* is a chain with empty boundary,  $\partial C = 0$ . A *boundary cycle* is a  $k$ -cycle that is the boundary of a chain in dimension  $k+1$ .

We use these definitions to construct various groups on a simplicial complex  $K$ . The set of all chains of dimension  $k$  forms the chain group  $C_k$ . The set of all cycles of dimension  $k$  forms the cycle group  $Z_k$ . The set of all boundary cycles of dimension  $k$  forms the a group  $B_k$ . The latter two groups can be understood in terms of the boundary operator  $\partial$  acting on  $K$ . The group  $Z_k$  is the kernel of the boundary operator,  $Z_k = \ker \partial_k$ . That is, it is the set of all  $k$ -chains that are sent to 0 by the boundary operator. The group  $B_k$  is the image of the boundary operator,  $B_k = \text{im } \partial_{k+1}$ . That is, it is the set of all  $k$ -chains which are themselves the boundary of  $(k+1)$ -chains in  $K$ . These groups have a particularly simple relationship to one another which is shown in Figure 2.11.

### 2.2.1.2 Homology

We are now ready to define homology, which will allow us to discuss and compare shape in a quantitative way. The  $j$ -th homology of a simplicial complex  $K$  is defined as the quotient group

$$H_j(K) = Z_j / B_j = \ker \partial_j / \text{im } \partial_{j+1}. \quad (2.7)$$

In words, homology is the group generated by equivalence classes of the cycle group  $Z_j$ , where equivalence is defined up to  $B_j$ . Elements of the homology group are classes of homologous cycles. Two  $j$ -cycles are homologous if they differ by the boundary of a  $(j+1)$ -chain. We work through a simple example in Figure 2.12.

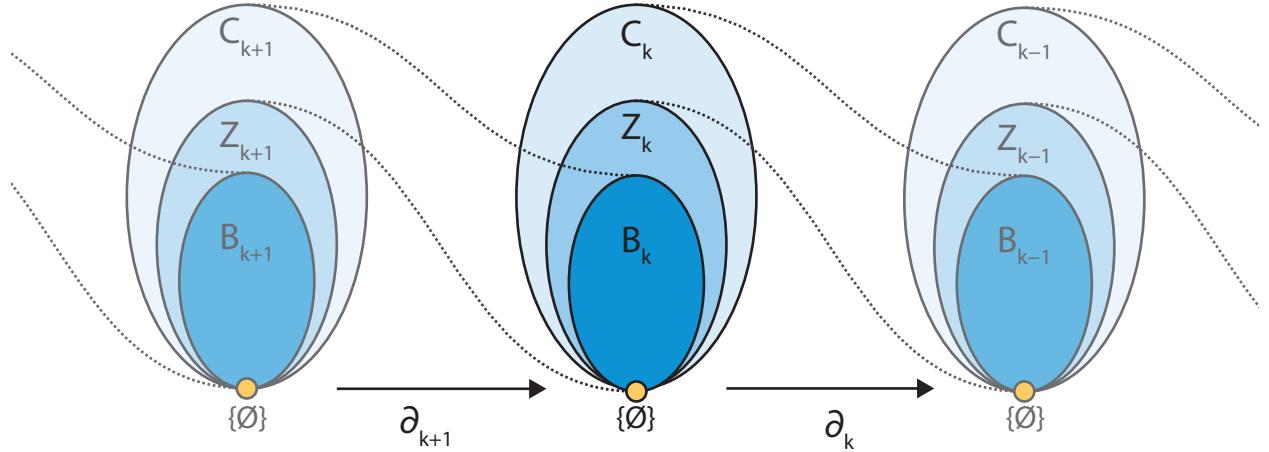


Figure 2.11: Relationship between the chain group ( $C_k$ ), cycle group ( $Z_k$ ), and boundary group ( $B_k$ ). Specifically,  $B_k \subset Z_k \subset C_k$ . We show the action of the boundary map  $\partial_k$  on each group at each dimension. Of particular note are the relations  $\partial_k : C_k \rightarrow B_{k-1}$  and  $\partial_k : Z_k \rightarrow \emptyset$ . Figure adapted from [62].

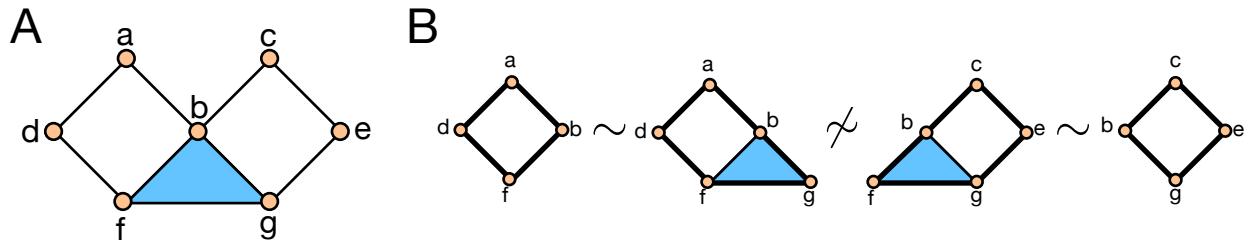


Figure 2.12: (A) A simplicial complex defined on a set of 7 vertices,  $S = \{a, \dots, g\}$ . The object has one connected component ( $\beta_0 = 1$ ) and two holes ( $\beta_1 = 2$ ). (B) Four cycles that can be defined on the complex. Cycles  $z_1 = \{(a, b) + (b, f) + (f, d) + (d, a)\}$  and  $z_2 = \{(a, b) + (b, g) + (g, f) + (d, f) + (d, a)\}$  are homologous, differing only by the cycle  $c_1 = \{(b, g) + (g, f) + (f, b)\}$  which is the itself the boundary of the closed triangle  $(b, f, g)$ . Likewise with cycles  $z_3$  and  $z_4$ . The two sets of cycles are not homologous with each other, and there are therefore constitute two independent elements of the homology group  $H_1(S)$ .

The rank of the homology group  $\|H_j(K)\|$  is the Betti number  $\beta_j$ . Intuitively, the Betti number represents the number of  $j$ -dimensional holes in the simplicial complex.

### 2.2.1.3 Constructing Complexes From Data

Finally, we must consider how to construct a simplicial complex from a given point cloud  $S$ .<sup>15</sup>

There are two common constructions we will describe: the *Čech complex* and the *Vietoris-Rips complex*. Both constructions involve a scale parameter  $\epsilon$ , and balls of radius  $\epsilon$  placed at the center of each vertex in  $S$ . Edges are drawn between vertices when balls overlap, that is, when  $d(v_a, v_b) < 2\epsilon$ . Where the two constructions differ is in how higher-dimensional simplices are filled in.

The Čech complex consists of the set of simplices  $\sigma$  with vertices  $s_1, \dots, s_k \in S$  such that

$$\text{Čech}(S, \epsilon) = \{\sigma \in S \mid \cap_i B(s_i, \epsilon) \neq \emptyset\}. \quad (2.8)$$

That is, the simplex  $\sigma_{(s_x, \dots, s_z)}$  is present if the intersection of balls of radius  $\epsilon$  centered on vertices  $(s_x, \dots, s_z)$  is nonempty. The Vietoris-Rips complex,  $VR(S, \epsilon)$ , is defined as

$$VR(S, \epsilon) = \{\sigma \in S \mid \text{diam}(\sigma) \leq 2\epsilon\} \quad (2.9)$$

where  $\text{diam}(\sigma) = \{\sup d(i, j) \mid i, j \in \sigma\}$ . In the Vietoris-Rips complex, a higher-dimensional simplex is filled in if every pairwise distance is less than  $2\epsilon$ . The difference between the two constructions is shown in Figure 2.13. In general,  $\text{Čech}(S, \epsilon) \in VR(S, \epsilon)$ .

The Čech complex is theoretically preferable because it comes with a *nerve theorem*, which states that the topology of the resulting complex will be equivalent to the topology of the union of balls used to create it. However, the Čech complex has drawbacks that prevent it from being widely applied to arbitrary data. First, computing the intersection of arbitrary balls is an expensive operation. While efficient algorithms exist in Euclidean space (the miniball algorithm [68]), it is much more difficult in arbitrary metric spaces. Furthermore, the Čech construction explicitly requires an ambient space in which the data is embedded. For data which comes in the form of a finite metric space, it may not be clear what is the

---

<sup>15</sup>In the author's view, this is the most important step in applying a TDA pipeline.

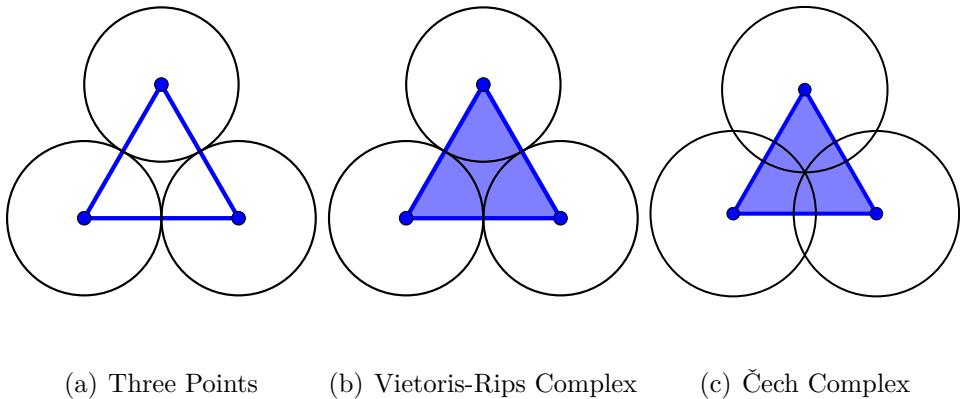


Figure 2.13: An example of the difference between Vietoris-Rips and Čech complex on an equilateral triangle. Consider each point to be 1 unit apart. In the Vietoris-Rips complex, the triangle is filled in when every pairwise edge is connected ( $\epsilon = 0.5$ ). In the Čech complex, the triangle is only filled in when all three balls intersect ( $\epsilon = 0.577$ ).

ambient space.<sup>16</sup> In practice, the Vietoris-Rips complex is more widely applied, because it requires only the set of pairwise distances between each vertex and a scale parameter  $\epsilon$ . The complex can be directly read off from the set of edges (known as the *1-skeleton*), making it extremely fast to compute.

### 2.2.2 Persistent Homology

Persistent homology is a tool developed under the umbrella of TDA that allows the shape of a point cloud to be computed across multiple scales simultaneously. Shape is quantified in terms of topological invariants representing homology, as discussed in the previous section. To understand why multiscale information might be of interest, consider the example in Figure 2.14. The data is sparse and noisy, but, to the eye, immediately appears to consist of two circles joined at a point along their edges. The two circles, however, are of a different radius. In the Figure, we show Vietoris-Rips complexes constructed at different scale parameters on the data. We observe that while some scale parameters are sufficient to resolve either one or

---

<sup>16</sup>This is indeed the case for the genomic data we consider: it is not immediately obvious what the intersection of three sequences defined over a finite alphabet should be. We discuss this further in Chapter 3.

the other of the two circles, no single scale parameter is sufficient to simultaneously capture the two shapes. Persistent homology solves this by providing an efficient way to track the shape information across *all* scale parameters.

Our object of study is the nested set of simplicial complexes, called a *filtration*, that is produced by tuning the scale parameter up to some threshold. At the smallest scale,  $\epsilon = 0$ , the complex consists only of disconnected points. As the scale parameter is increased, the topology of the complex changes – clusters merge, holes and loops form, other holes and loops are filled – until the complex is fully connected. Each aspect of shape represents a topological invariant, and as the scale is changed, the birth and death of different invariants is encoded as an interval  $(b_i, d_i)$ .

The shape information can be concisely summarized in a *barcode diagram*. The barcode diagram represents topological features as horizontal line segments, annotated with a birth-death interval, and indexed by dimension. The birth time is when a particular invariant first appears in the complex, and the death time is when the invariant collapses in the complex.  $H_0$  represents the number of connected components and is roughly equivalent to a hierarchical clustering of the data. Higher dimensions represent loops ( $H_1$ ), voids ( $H_2$ ), and their generalizations in the data. The number of bars at a particular scale will be the Betti number  $\beta_n(\epsilon)$  for the complex  $K(\epsilon)$ . Taken together, the barcode diagram represents a complete and quantitative picture of the shape of the data.

The information can be equivalently represented as a persistence diagram, which is a scatter plot of invariants with birth time on the  $x$ -axis and death time on the  $y$ -axis. The barcode diagram and persistence diagram for the two circles data is shown in Figure 2.15. First, looking at  $H_0$ , we see that the data begins disconnected and becomes connected at around  $\epsilon = 24$ . Next, looking at  $H_1$ , we count eight loops across a range from  $\sim 5$  to  $\sim 80$ . Two of these loops persist for what appears to be an appreciable length of time. We associate these two loops with the two circles that we identified qualitatively from the raw point cloud data.

The intuition behind persistent homology is exactly that: good or interesting topological features will be robust and persist over long scales. In the barcode diagram, this corresponds to longer bars; in the persistence diagram, this corresponds to points sitting far from the diagonal. Invariants that we observe persisting for only short scales are likely to be noise or other artifacts.<sup>17</sup> And because a single scale is not capable of representing all features of the data, we examine all scales simultaneously.

In fact, the persistence algorithm is more powerful than that, and can return not only the intervals associated with the invariants, but *representative cycles* of each invariant. The representative cycles correspond to a set of simplices that surround an invariant, and can be used to determine which data points are involved in a particular invariant.

To summarize, a complete description of the persistent homology pipeline is shown in Figure 2.16. The pipeline is as follows: A dataset,  $S = (s_1, \dots, s_N)$ , is represented as a point cloud in a high-dimensional space (not necessarily Euclidean). From the point cloud, a nested series of simplicial complexes, or a filtration, is constructed, parameterized by a filtration value  $\epsilon$ . The filtration is represented as a list of simplices defined on the vertices of  $S$ , annotated with the  $\epsilon$  at which the simplex appears. Given a filtration, the persistence algorithm is used to compute homology groups. The 0-dimensional homology ( $H_0$ ) represents a hierarchical clustering of the data. Higher dimensional homology groups represent loops, holes, and higher dimensional voids in the data. Each feature is annotated with an interval  $(b_i, d_i)$ , representing the  $\epsilon$  at which the feature appears and the  $\epsilon$  at which the feature collapses in the filtration. These filtration values are the *birth* and *death* times, respectively. The set of intervals are represented as either a barcode diagram or a persistence diagram.

A second way of applying persistent homology is through sublevel sets. The sublevel set of a function  $f : X \rightarrow \mathbb{R}$  is defined as

$$L_\epsilon^-(f) = \{x \mid f(x) \leq \epsilon\}. \quad (2.10)$$

---

<sup>17</sup>The obvious question of how to rigorously determine what makes a good interval is an open question that is currently being addressed by a number of different groups. We discuss this further in Section 2.2.2.2.

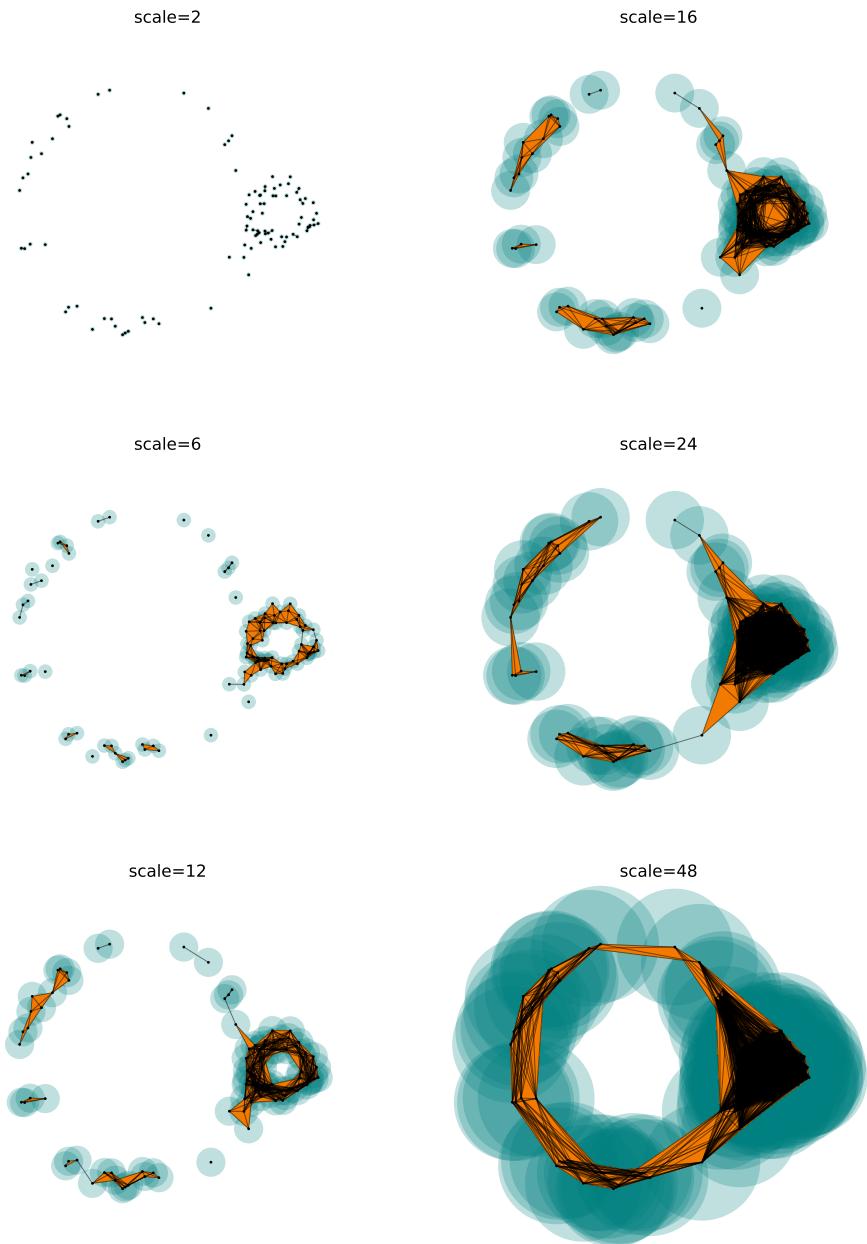


Figure 2.14: An example of constructing a filtration. The nested series of complexes form a filtration. Persistent homology will compute and track the homology at each scale.

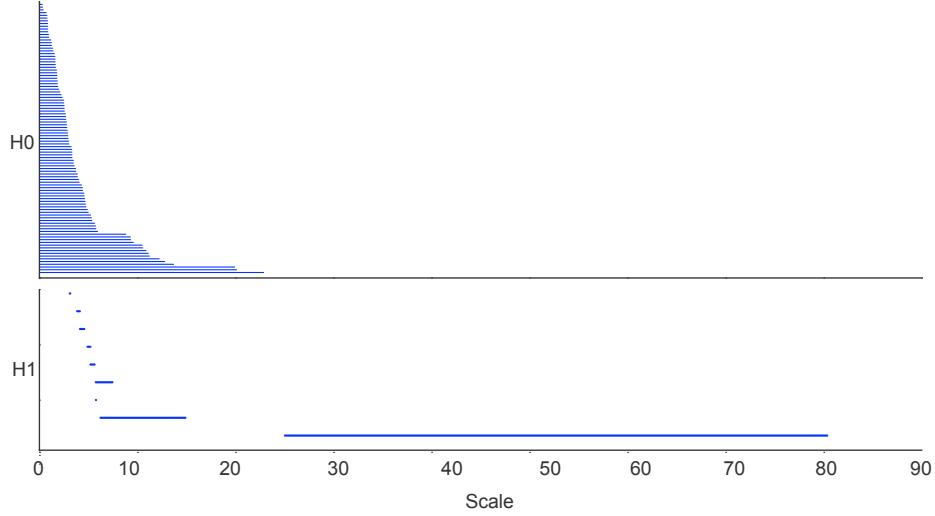


Figure 2.15: Barcode Diagram for the example in Figure 2.14.  $H_0$  and  $H_1$  represent connectivity and holes, respectively. The two holes visually apparent in the data are reflected in the two long bars in the diagram, persisting across different scales. Shorter bars are interpreted as topological noise.

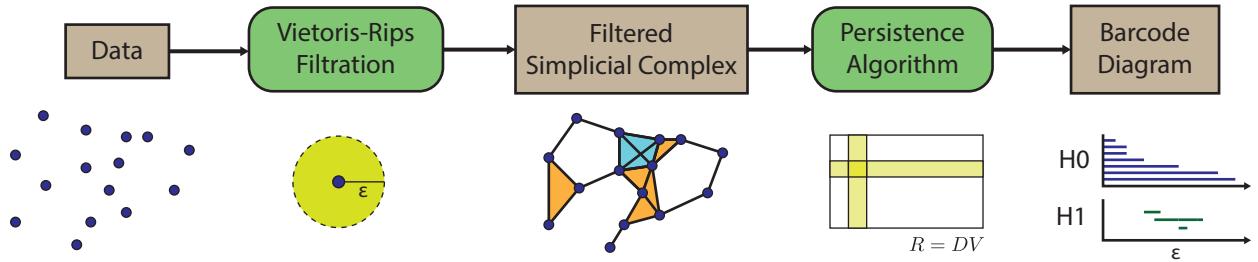


Figure 2.16: The Persistence Pipeline. Data is represented as a high-dimensional point cloud. A Vietoris-Rips filtration, a nested series of simplicial complexes parameterized by a scale  $\epsilon$  is constructed. Given a filtration, the persistence algorithm is used to compute homology groups. Topological features are represented by an interval  $(b_i, d_i)$ , representing the birth and death times of the feature. The set of intervals are represented as either a barcode diagram or a persistence diagram.

One can compute persistent homology of sublevel sets by varying the parameter  $\epsilon$ . In this case, the beginning and end of each bar will correspond to critical points of the function. A simple example is shown in Figure 2.17.

As primarily end-users of persistent homology, the details of the persistence algorithm are largely beyond the scope of this thesis. Effectively, it involves manipulating the boundary

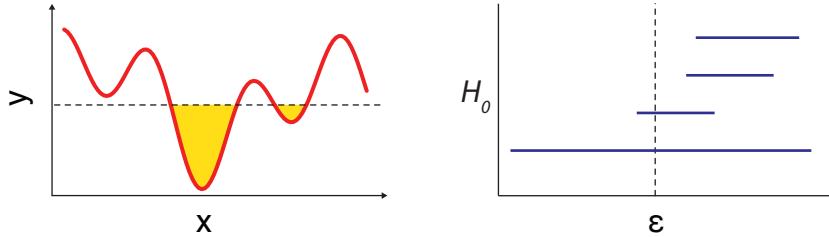


Figure 2.17: In a sublevel-set filtration, the homology of sublevel-sets of a function  $f : \rightarrow \mathbb{R}$  are computed. Here we illustrate this for a simple function, showing the sublevel set at a particular  $\epsilon$  and its position in the barcode diagram indicating two connected components. In general, the beginning and end of bars will correspond to critical values of  $f$ .

matrix into a particular reduced form, from which each bar and representative cycle can be read off. Several packages for computing persistent homology have been developed, including Javaplex [142], Dionysus [144], Perseus [117], Gudhi [111], and PHAT [13]. Additionally, the R TDA package wraps functions from Dionysus and Gudhi in a user-friendly frontend [63].<sup>18</sup>

### 2.2.2.1 Stability of the Persistence Algorithm

The stability statement gives a foundation for comparing the persistence diagrams from different data. Of particular interest, it gives conditions on the effect of noise on data sampled from a particular object. The stability result guarantees that small perturbations in the input data will produce only small changes in the output diagrams. The result is due originally to Chazal *et al.* [32]. The statement requires two things: (1) a notion of distance between the input data  $D$  and the perturbed data  $D'$ , and (2) a notion of distance between the resulting diagrams  $B$  and  $B'$ .

First, we need a notion of distance between persistence diagrams. Recall the persistence diagram consists of a set of intervals  $(b_i, d_i)$  along with the diagonal. We introduce the concept of a *matching*. For two persistence diagrams  $B$  and  $B'$ , a matching is simply a

---

<sup>18</sup>In our work we have relied on a variety of these packages. For straight-forward construction of the barcode diagram, we find the R package TDA easiest to use. If one needs to directly build and manipulate filtered simplicial complexes, Dionysus has convenient Python bindings. For large datasets, PHAT and its parallel implementation DIPHA are recommended [11, 12].

mapping of intervals in  $B$  to intervals in  $B'$ , where we allow points to match to the diagonal (to account for cases with an unequal number of points). For a matched pair of intervals  $(a, b)$ , we define the  $L_\infty$  distance as

$$d_\infty(a, b) = \max\{|a_x - b_x|, |a_y - b_y|\}. \quad (2.11)$$

The *bottleneck cost* of a matching between two diagrams is defined as the maximum  $L_\infty$  distance among all matched pairs. The *bottleneck distance* is defined to be the minimal bottleneck cost across all matchings,

$$d_B(B, B') = \inf_{\gamma \in \Gamma} \sup_{p \in B} \|p - \gamma(p)\|_\infty \quad (2.12)$$

The matching with minimal bottleneck cost is the *bottleneck matching*.

Second, we need a notion of distance between finite metric spaces. Here we will use the *Gromov-Hausdorff distance*, which measures how far apart two spaces are from being isometric. It measures the longest distance from a point in one set to the closest point in another set within a metric space.

$$d_{GH}(X, Y) = \inf_{f, s} d_H(X, Y) \quad (2.13)$$

The result of [32] states that the bottleneck distance between  $B$  and  $B'$  is bounded by the Gromov-Hausdorff distance between the finite metric spaces embedded in  $A$  and  $B$ .

$$d_B(H_K(X), H_K(Y)) \leq d_{GH}(X, Y) \quad (2.14)$$

The idea is easiest to visualize using a level-set persistence example, which we show in Figure 2.18. Here we see a simple function,  $f(x)$ , and a noisy sample of the same function. The persistence diagrams of each are shown in Figure 2.18B, along with a bottleneck matching, shown in yellow. As can be seen, most of the noise introduced is reflected in intervals close to the diagonal. While this example was for a simple level-set filtration, the result has an extension to the arbitrary finite metric spaces  $(X, d_X)$  which we primarily consider in this thesis.

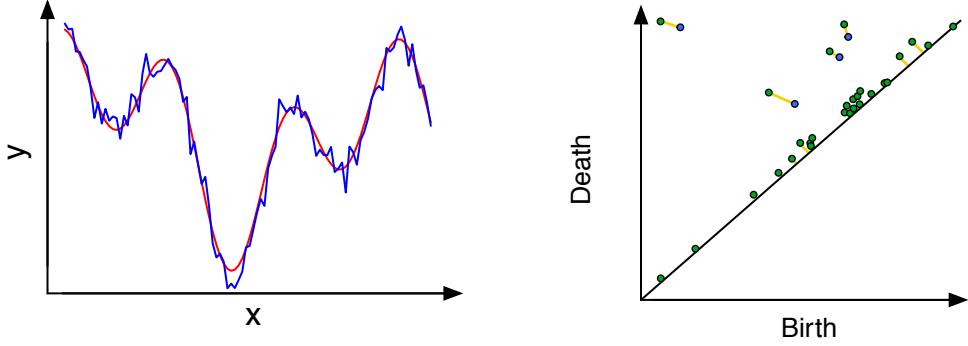


Figure 2.18: An example of the stability of the persistence diagram with respect to noise. (A) A function (red), and a noisy sampling of the same function (blue). (B) The level-set persistence diagrams for the two samples. A bottleneck matching is indicated with yellow lines. The stability result gives an upper bound on how much the diagrams can diverge as the sample diverges from the true function. This figure is adapted from [56].

### 2.2.2.2 Statistical Persistent Homology

Persistent homology has been developed largely as an exploratory tool for data analysis. However, one would like to integrate it in analysis pipelines more broadly, which requires notions of statistics to be developed. One of the difficulties is that the persistence diagram, consisting of a multiset of points in the plane, can be somewhat unwieldy to work with. Substantial recent work in the TDA community has focused on these questions in order to develop statistical foundations for persistent homology. We give here a brief flavor of some of these ideas and their relation to our own work. There are three main threads in statistical persistent homology:

1. Confidence intervals defined on the persistence diagrams
2. Functional summaries of the persistence diagram
3. Probability measures on the space of persistence diagrams

Confidence intervals on the persistence diagram address the question of when a bar is a significant topological feature and when it should be considered topological noise. Fasy *et al.* have developed ways of generating confidence intervals for persistence diagrams [62]. They use a filtration on a kernel density estimate of the data, and bootstrap resampling,

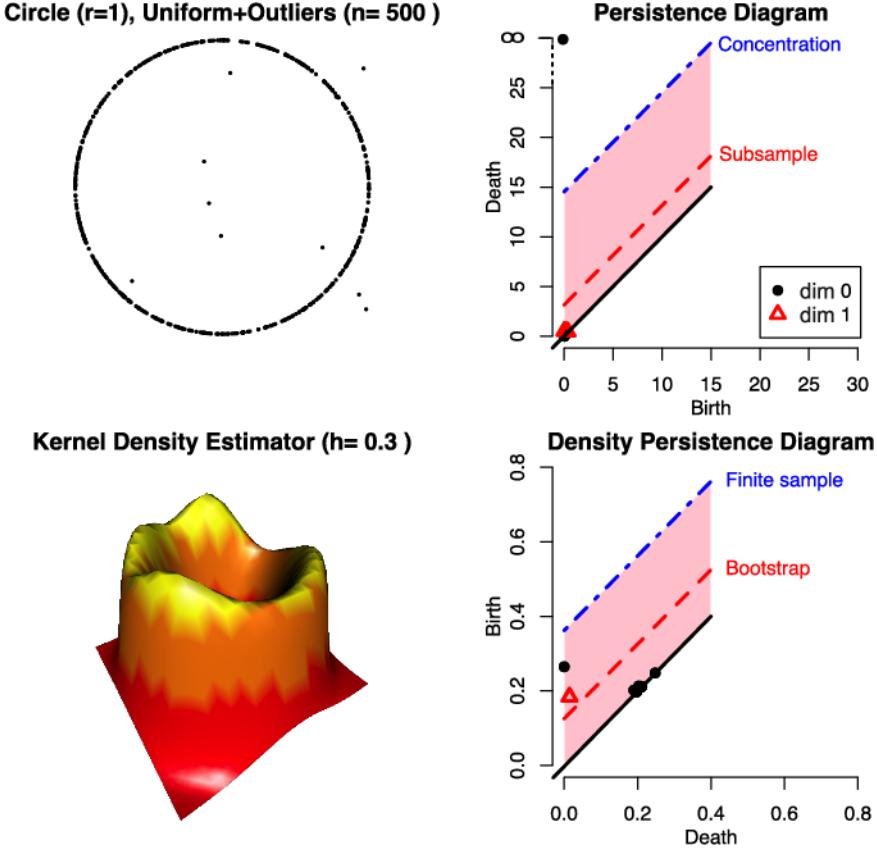


Figure 2.19: An example of how statistical persistent homology can be used both to handle noise in data and to put confidence sets on the persistence diagram. Figure is taken from [62].

to place an off-diagonal line on the diagram, below which points are to be considered noise. This approach can be used to handle outliers and noise in the data by providing a way to identify the most robust topological features. An example of estimating the topology of a circle, with outliers, is shown in Figure 2.19. Using the density estimate with a bootstrap estimator one can recover the circle as a significant feature. Subsampling estimates were studied further in [33], and related approaches were developed by Blumberg *et al.* in [16].

Functional summaries of the persistence diagram convert the multiset of intervals into . From functional summaries, machine learning approaches can be used downstream. Bubenik has developed the language of persistence landscapes [24, 23]. Essentially, the landscape is generated by rotating the persistence diagram 45 degrees and dropping a tent at each point.

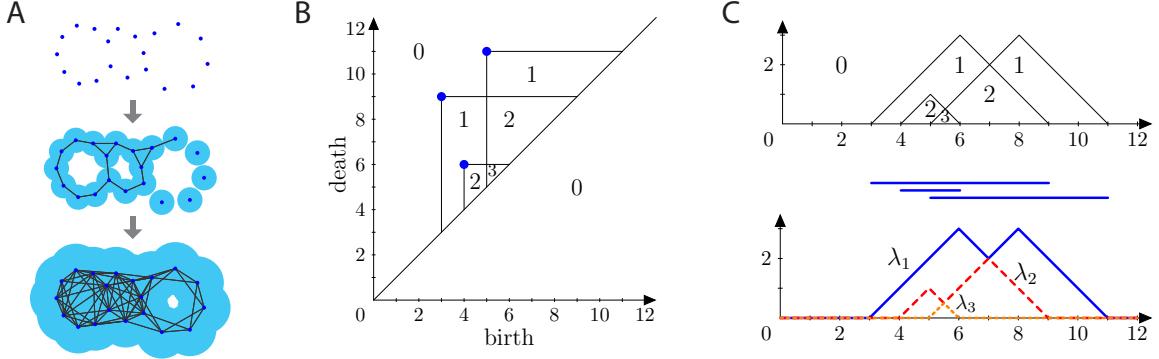


Figure 2.20: Persistence Landscapes. (A) A simple filtration on a set of points. (B) The persistence diagram from this data. (C) Transformed into a persistence landscape. This figure is adapted from [23, Fig. 2].

The silhouette is taken by weighting the contribution of each point based on its height.

$$\Lambda_p(t) = \begin{cases} t - x + y & t \in [x - y, x] \\ x + y - t & t \in (x, x + y] \\ 0 & \text{otherwise} \end{cases} = \begin{cases} t - b & t \in [b, \frac{b+d}{2}] \\ d - t & t \in (\frac{b+d}{2}, d] \\ 0 & \text{otherwise.} \end{cases} \quad (2.15)$$

An example of the persistence landscape and silhouette is shown in Figure 2.20.

A second approach has been taken by Schweinhart and MacPherson in [109], in which a transformation of the intervals is used to characterize random polymers in terms of a fractal dimension. A related approach has recently been proposed by Kwitt *et al.* that represents the persistence diagram in a kernel space for use in machine learning [98, 129]. We explore the use of functional summaries of the barcode diagram in a statistical inference setting in Chapter 4.

Finally, probability measures on the space of persistence diagrams. Several authors have examined the space of persistence diagrams as a Polish space, with well-defined notions of mean and variance. See the work of Turner [146] and Mileyko [113]. This work crucially relies on distances between diagrams that are based on matchings, as discussed in Section 2.2.2.1. Further work has focused on statistical properties of the persistence diagrams themselves, such as [34]. Establishing these foundations would lead to direct use of the persistence diagram

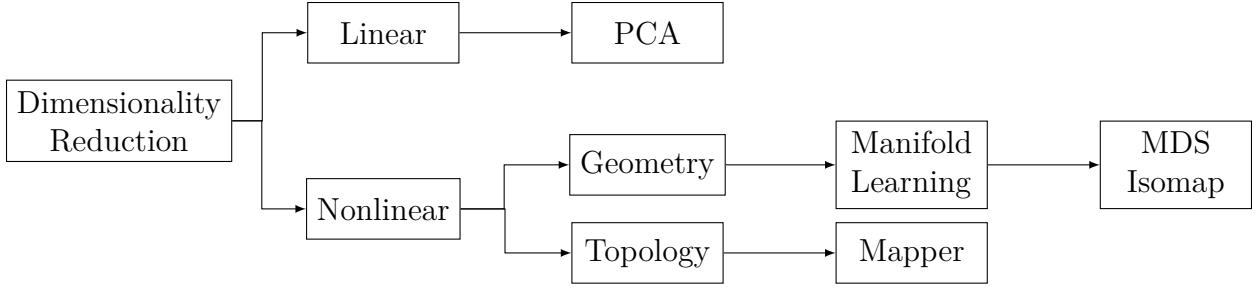


Figure 2.21: Dimensionality Reduction Approaches for Exploratory Data Analysis.

as a statistical object, but as of yet no efficient tools have been developed.

### 2.2.3 Mapper

Mapper is an approach for the representation and visualization of patterns in high-dimensional data. As such, it sits within the larger class of dimensionality reduction algorithms for exploratory data analysis (EDA), such as multidimensional scaling (MDS) [97], Isomap [143], and t-SNE [147]. A perspective on how these various approaches to EDA are related is shown in Figure 2.21. The primary distinction between Mapper and the existing class of nonlinear dimensionality reduction algorithms is that Mapper seeks to preserve the topology of the input data, rather than geometry. Compared to existing approaches, Mapper has the following advantages: (1) it is coordinate free, depending only on the metric properties of the data; (2) an invariance to deformation, which provides a robustness against noise; and (3) results in a compressed representation, which gives the ability to handle extremely large data. Mapper was initially developed by Singh and Carlsson in [134]; further exposition and examples can be found in [107]. Mapper has been applied to problems in RNA folding [20], breast cancer subtype classification [122], and genetic associations in type 2 diabetes [103].

A simple example of the Mapper algorithm is shown in Figure 2.22, which is adapted from [107]. As input is a point cloud  $X$  with an associated metric, Euclidean or not (Figure 2.22)A. First, a filter function is applied to the data (Figure 2.22B). The filter function maps the original points onto the real line,  $X \rightarrow \mathbb{R}$ . Standard filters include things like the

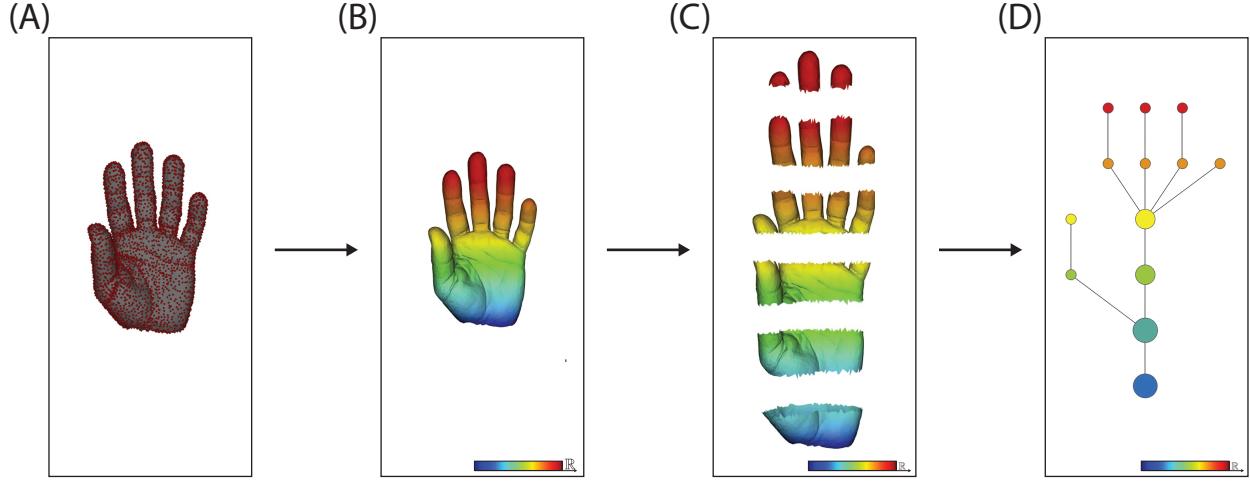


Figure 2.22: The Mapper Algorithm. Mapper starts with a set of data points and a filter function  $f$ , and produces as output a graph that captures the shape of the data. (A) The original data represented as a point cloud. (B) The filter function  $f$  projected onto the data. (C) In the projected space, the data is divided into overlapping bins. (D) Individual bins are represented as nodes. The size of each node represents the number of data point contained in the bin. Nodes can be colored by various attributes by the data, in this example the average value of the filter function is used. Pairs of bins that have points in common are connected by an edge.

mean, density,  $L_1$ -centrality, and the first and second components of a PCA decomposition. Second, the projected space is split into overlapping bins based on a resolution (bin size) and overlap parameter.<sup>19</sup> Third, the bins are then clustered in the original high-dimensional space (Figure 2.22). Each cluster will form a node in the graph representation.<sup>20</sup> Finally, nodes that share points in the original space are connected by an edge.

Our use of Mapper will be primarily as a means of visualizing relationships in sequence data. While the resulting graphs cannot be strictly interpreted in a phylogenetic sense, they will provide valuable information about evolutionary relationships. We use the commercial implementation of Mapper developed by Ayasdi [6]. Open-source implementations of Mapper are available in the Python Mapper package [115] and sakmapper package [160].

<sup>19</sup>Multiple filter functions can be used and binned on a grid.

<sup>20</sup>Nodes in a Mapper graph consist of multiple points in the original data; this is essence of the compressed representation.

## 2.3 Applying TDA to Molecular Sequence Data

Aligned molecular sequence data can be naturally viewed as a point cloud in a high-dimensional space, which we loosely call *sequence space*. The particular structure of sequence space will be determined by the length,  $L$ , of the aligned sequences, and the alphabet,  $Q$ , over which the sequences are defined. The typical sequence alphabet will be either nucleotides or amino acids. The dimension of the space is determined by  $L$ . Sequence space will therefore consist of the  $||Q||^L$  possible sequences. Together with any of the standard genetic distance measures, this forms a metric space.

The processes of evolution can be seen as an exploration of sequence space. Clonal evolution, in which mutations accrue at each generation, will be Clonal evolution is the process of smoothly moving through sequence space, while reticulate evolution is the process of making discontinuous jumps through the space. Our data consists of a subset of points sampled from sequence space. These points reflect a particular evolutionary history. As more data is acquired, regions of sequence space will become more densely sampled and our ability to reconstruct evolutionary relationships will be improved.

Given sequence data, our program is to (1) encode the data as a finite metric space, (2) use tools from TDA to characterize the topology of the data, and (3) interpret the topology in an evolutionary context. In particular, we apply persistent homology, and read phylogenetic information contained in the dataset off the resulting barcode diagram. This idea was first proposed in [31]. In that paper, the authors developed two metrics for measuring reticulate evolution using homological features: (1) topological obstruction to phylogeny (TOP), which uses the  $L_\infty$ -norm of the barcode as a coarse measure of reticulation; and (2) irreducible cycle rate (ICR), which uses temporal annotations to measure the average number of  $H_1$  features per unit time. They applied this approach to a variety of viral datasets, including influenza and HIV. The work presented in this thesis extends this work in several substantial ways. We make two preliminary remarks before considering a more complex example.

### 2.3.1 Topology of Tree-like Metrics

An important foundational point was demonstrated by Carlsson in [31]. Recall that tree-like data will have an additive metric, as described in Section 2.1.4.2. In [31], it was proven that for additive metric spaces, the Vietoris-Rips filtration of the data will consist of a nested set of acyclic complexes. Consequently, the persistent homology of additive data will have nontrivial topology only in dimension zero. Furthermore, while noise in the data will introduce small deviations from additivity, the theorem puts bounds on the size of the topological features that can arise in this manner. These bounds rely on the Gromov-Hausdorff stability conditions described in Section 2.2.2.1. On the other hand, if the evolutionary history includes reticulate events that cannot be represented as a tree, these events will be captured as non-trivial higher dimensional homology in the barcode diagram, an idea which we develop below. This theorem provides an important negative control in using TDA to characterize reticulate evolution.

### 2.3.2 The Fundamental Unit of Reticulation

In population genetics, there is a simple test for the presence of reticulate evolution in sequence data called the *four-gamete test* [82]. The test assumes only an infinite-sites model, which states that for a sufficiently long genome, a particular residue can only ever undergo a single mutation. Put another way, there is no multiple-mutation or back mutation. The infinite-sites model has three consequences: first, one need only consider segregating sites, or nucleotide positions that have undergone a mutation. Second, because a given position can mutate only once, it is sufficient to represent sequences as binary strings, where a 0 indicates the unmutated state and 1 the mutated state. Third, for a given position we can arbitrarily assign the unmutated and mutated states. The infinite-sites model is considered a reasonably good model for long genomes.

The four-gamete test identifies reticulate evolution by looking at pairs of segregating sites. Given biallelic data, there are four possible haplotype patterns, or states, for a pair

of segregating sites: 00, 10, 01, or 11.<sup>21</sup> The statement of the four gamete test is this: in any given dataset, the simultaneous presence of all four haplotype states in any pair of segregating sites is incompatible with strictly clonal evolution, and indicates reticulate evolution. To see this, assume state 00 as the ancestor to states 10 and 01, which arise from two independent mutations. Because of the no multiple-mutation assumption, it is not possible for either of these two states to then independently mutate into state 11. The only way for state 11 to arise is via a reticulate event that brings together the left site from state 10 and the right site from 01.<sup>22</sup> This process is illustrated in Figure 2.23A.

Under a Hamming metric, the distance matrix for the set of four sequences  $s_1 = 00$ ,  $s_2 = 10$ ,  $s_3 = 01$ , and  $s_4 = 11$  is

$$d = \begin{pmatrix} 0 & 1 & 1 & 2 \\ 1 & 0 & 2 & 1 \\ 1 & 2 & 0 & 1 \\ 2 & 1 & 1 & 2 \end{pmatrix} \quad (2.16)$$

The Vietoris-Rips filtration of this space is shown in Figure 2.23B. At  $\epsilon = 0$  the four sequences are disconnected. At  $\epsilon = 1$ , four edges are drawn, forming a loop. At  $\epsilon = 2$ , the space is completely connected and the loop is killed. Persistent homology captures the presence of this loop as an  $H_1$  feature in the interval [1, 2) (Figure 2.23C). In this way, the reticulate event is associated with the presence of a nonzero  $H_1$  bar. A possible reticulate evolutionary genealogy representing this data, including two mutations ( $m_1$  and  $m_2$ ) and one reticulation ( $r_1$ ) is shown in Figure 2.23D.

We consider this example to be the minimal, or fundamental, unit of reticulation. More complicated patterns of reticulation can be seen as extensions of this example.

<sup>21</sup>These sites need not be adjacent.

<sup>22</sup>It is entirely possible for the reticulate event to have had a reversed pattern of ancestry, in which case the reticulation would result in a state 00 and would not be detectable from the sequence data.

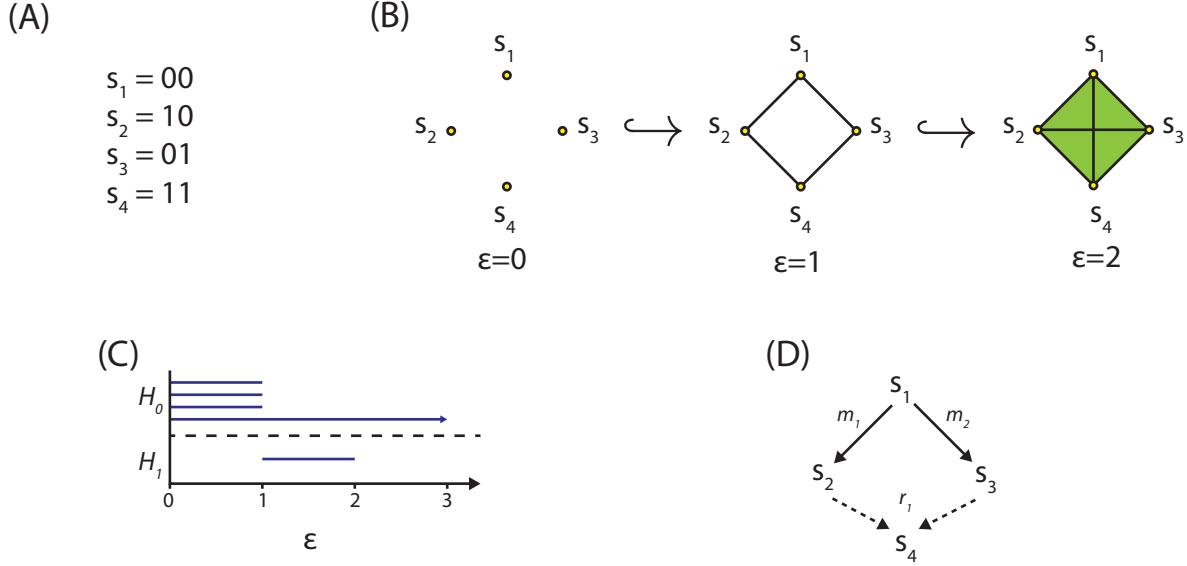


Figure 2.23: The Fundamental Unit of Reticulation. (A) A set of four sequences. (B) Vietoris-Rips filtration of the sequences. (C) Barcode diagram of the filtration. (D) An evolutionary genealogy including two mutations ( $m_1$  and  $m_2$ ) and a single reticulation ( $r_1$ ).

### 2.3.3 A Complete Example

We illustrate a complete example of how TDA can capture reticulate evolution from complex population data in Figure 2.24. Consider the reticulate phylogeny (Figure 2.24A): five genetic sequences sampled today (yellow circles) originate from a single common ancestor due to clonal evolution (solid blue lines tracing parent to offspring) and reticulate evolution (dotted red lines). In Figure 2.24B, these five samples are placed in the context of a larger dataset, where the data has been projected onto the plane using PCA. Persistent homology is then applied to this larger sample. In Figure 2.24C we demonstrate the construction of a filtered simplicial complex, showing how the connectivity changes as the scale parameter  $\epsilon$  is increased. Finally, in Figure 2.24D we see the resulting barcode diagram. Using  $H_0$  we can track the number of strains or subclades that persist, roughly corresponding to the tree-like component of the data. The  $H_1$  bar spanning roughly  $\epsilon = 0.13$  to  $\epsilon = 0.16$  identifies the presence of a reticulate event involving the five highlighted sequences. The scale over which this bar persists represents the amount of evolutionary time separating the parents and the

Table 2.2: Dictionary connecting algebraic topology and evolutionary biology

Algebraic Topology	Evolutionary Biology
Filtration value $\epsilon$	Genetic distance (evolutionary scale)
0-dimensional Betti number at filtration value $\epsilon$	Number of clusters at scale $\epsilon$
Generators of 0-D homology	A representative element of the cluster
Hierarchical relationship among generators of 0-D homology	Hierarchical clustering
1-D Betti number	Lower bound on number of reticulate events
Generators of 1-D Homology	Reticulate events
Generators of 2-D Homology	Complex horizontal genomic exchange
Non-zero high-dimensional homology (topological obstruction to phylogeny)	No treelike phylogenetic representation exists
Number of higher-dimensional generators over a time interval (irreducible cycle rate)	Lower bound on recombination/reassortment rate

reticulate offspring. Additionally, the persistence algorithm will return a generating basis for a particular homology group, which we can use to identify the particular mixtures of sequences involved a reticulation. In this way, we can analyze both the scale and frequency of reticulation in genomic data sets.

We summarize the connection between genomic data and TDA in Table 2.2.

### 2.3.4 The Space of Trees, Revisited

In Section 2.1.4.4, tree space was introduced as an abstract construction to systematically represent the set of all possible binary trees as a geometric space. Tree space on  $n$  leaves,  $\mathcal{T}_n$ , was shown to be the subspace of the complete space of finite metrics on  $\mathbb{R}_{\geq 0}^{\binom{n}{2}}$  consisting of those metrics that satisfy the four-point condition (or additivity). Because real sequence data will very rarely satisfy this condition, one possible interpretation of phylogenetics is of finding the best projection onto tree space for arbitrary data.

The program we propose can be understood as an extension of the tree space framework. We would like to use topological invariants, specifically homology, to measure the frequency and scale of reticulate evolution in sequence data. First, from the theorem due to Carlsson,

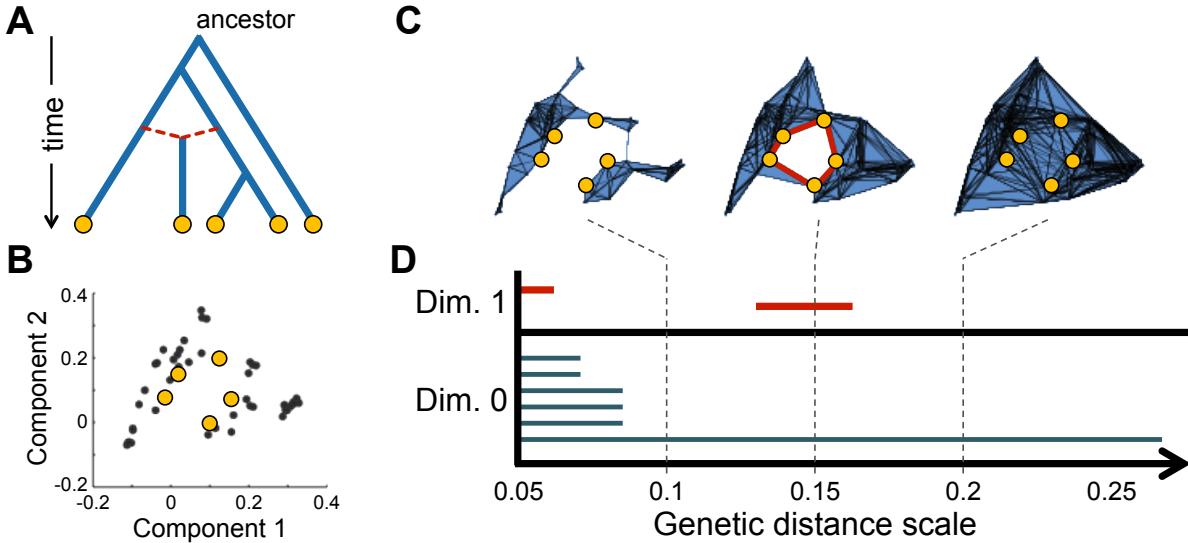


Figure 2.24: Applying persistent homology to genomic data. (A) An evolutionary genealogy including reticulation. (B) Data projected into 2-dimensions. (C) Construction of a filtered simplicial complex. (D) The resulting multiscale barcode diagram.

we know that higher homology will vanish on tree space. Second, from the simple example in Section 2.3.2, we know that non-additive reticulate processes will have nonvanishing higher homology. Rather than attempt to characterize arbitrary data by projecting onto tree space, we will use persistent homology to compute homological invariants that will characterize our space. Our hypothesis is that as the data moves further from tree space, it becomes increasingly nonadditive, which can be captured quantitatively with increasing signal from higher homology. Our updated picture is shown in Figure 2.25, where we depict tree space embedded in the larger space of finite metrics. We anticipate for appreciable datasets, our sensitivity will be such that those close to the tree space will have little to no homological signal, while those further away will have an increasing homological signal. A second point is that the metric structure will now allow us to pass through regions of space that are nonadditive. Hence, one could conceivably compare two trees by drawing the direct path between them in metric space, and evaluating the persistent homology at each point. While some interesting work has explored the combinatorial structure of the space of metrics on low numbers of points (see [138]), it does not appear feasible in general to provide a complete

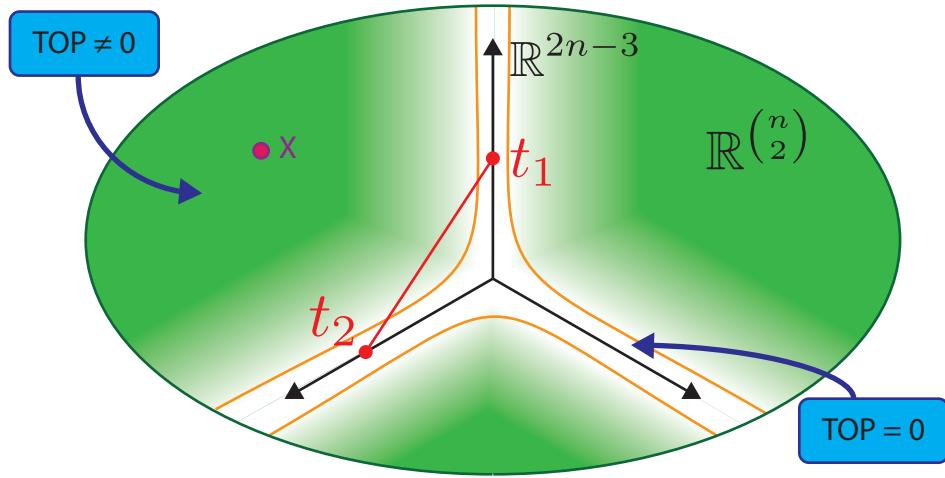


Figure 2.25: Tree space on  $n$  leaves,  $\mathcal{T}_n$ , is a subspace of the larger space of the metric cone on  $\mathbb{R}_{\geq 0}^{\binom{n}{2}}$ . In the presence of reticulate evolution, data may not sit near an additive tree, as for example data  $X$  (pink circle). The invariant TOP (topological obstruction to phylogeny) is one way to characterize these spaces and can be computed using persistent homology. As one moves further from tree space TOP increases. The orange lines indicate regions close to tree space in which TOP is insensitive to non-additivity. In this example we also indicate two trees  $t_1$  and  $t_2$  which sit on subspaces of different topology in tree space. The direct path between the two topologies will pass through a  $\text{TOP} \neq 0$  region.

decomposition.

# Part I

# Theory



# Chapter 3

## Quantifying Reticulation Using Topological Complex Constructions Beyond Vietoris-Rips

### 3.1 Introduction

In Chapter 2, the Vietoris-Rips complex was introduced as a construction on molecular sequence data. The persistent homology of a filtered sequence of complexes was shown to provide a quantitative measure of reticulate processes. As we will show, in certain cases the Vietoris-Rips complex can have a reduced sensitivity to reticulation. In this chapter we introduce two ideas to increase the usefulness of the signal generated by persistent homology. The first is an approach for imputing latent ancestors into the data that increases the quantitative signal detected by persistent homology. Our approach is built on the *median graph* construction. Median graphs form the basis for a large number of phylogenetic network algorithms and are closely related to split decompositions of finite metrics [8, 7]. A common desire is an approach to quantify the complexity of the resulting construction. We show that the persistent homology of the median closure is a fast and efficient way to identify the

phylogenetic incompatibility in a dataset. The second is an approach for computing Čech complexes from genomic data. The Čech complex, introduced in Section 2.2.1.3, has certain advantages over the Vietoris-Rips complex. However, it requires a notion of embedding space for data, which for genomic data is not entirely obvious.

The structure of this chapter is as follows. In Section 3.2 we show simple examples of the reduced sensitivity of the Vietoris-Rips for detecting reticulations. In Section 3.3 we introduce the median closure of the original vertex set. We show how this operation recovers invariant signals of phylogenetic incompatibility in a quantitative way. In Section 3.4 we present a Čech complex construction on sequence data. Throughout, we assume biallelic data under an infinite sites model with no back mutation.

## 3.2 Sensitivity of the Vietoris-Rips Construction

The fundamental loop (00,10,01,11) was introduced in Section 2.3.2 as the simplest example in which binary sequence data would manifest reticulation, as measured by persistent homology. The fundamental loop is based on the four-gamete test of haplotype incompatibility in an infinite sites model. In considering further small examples of sequence data we often encountered situations in which the four gamete test indicated reticulate evolution, but persistent homology failed to detect a loop. This was often due to degeneracies that would arise because of either (a) incomplete sampling in which case recombinations failed to be detected because parental and ancestral strains would collapse prior to connecting with the recombinant offspring, or (b) cases in which the recombination event led to an offspring that sat spatially intermediate to the ancestral and parental strains. We demonstrate with two examples.

**Example One** It is generally the case that we do not have a complete sampling of the sequences corresponding to the evolutionary history of a set of sequences. For example, we may not have sampled the true recombinant child, only a descendant which has accumulated

additional mutations. Consider the sequences  $s_1 = 000$ ,  $s_2 = 100$ ,  $s_3 = 010$ , and  $s_4 = 111$ . The four-gamete test identifies incompatibility between sites 1 and 2. However, persistent homology of the four sequences does not capture this reticulation. To understand why, consider  $s_1$  to be the common ancestor,  $s_2$  and  $s_3$  to be parents, and  $s_4$  to be a descendant of a reticulate event. In this scenario, we can infer that there was an ancestral recombinant sequence,  $s_r = 110$ , which was not sampled. The failure to find a loop is due to the ancestral and parent sequences collapsing before connecting with the recombinant offspring, as shown in Figure 3.1A. In general, for a loop to be detected, the two internal distances must be greater than any of the four external distances. In this case, the internal distance from parent 1 ( $s_2$ ) to parent 2 ( $s_3$ ),  $d_{23}$  is equal to the distances from each parent to the sampled descendant of the recombinant ( $d_{24}$  and  $d_{34}$ ). This is an example of incomplete sampling lowering the detection sensitivity, even in cases where incompatible sites are present.

**Example Two** This example is taken from [137]. Consider the sequences:  $s_1 = 0000$ ,  $s_2 = 1100$ ,  $s_3 = 0011$ ,  $s_4 = 1010$ , and  $s_5 = 1111$ . The four-gamete test identifies incompatibilities between sites 1 and 3, 1 and 4, 2 and 3, and 2 and 4. Performing the Hudson-Kaplan test yields a partition between sites 2 and 3, however [137] show a minimum of two reticulate events are required to explain the data. Using the standard filtration, the complex contracts completely at  $\epsilon = 2$ , and no higher homology will be detected. In this case, the two reticulations interact in such a way that  $s_3$  now sits equidistant from the other four sequences. Had  $s_3$  not been in the data, we would have had an example very similar to Example 1, with the interpretation of one recombination event. In this example we observe that multiple reticulate events can interact in complicated ways, obscuring the signal from persistent homology.

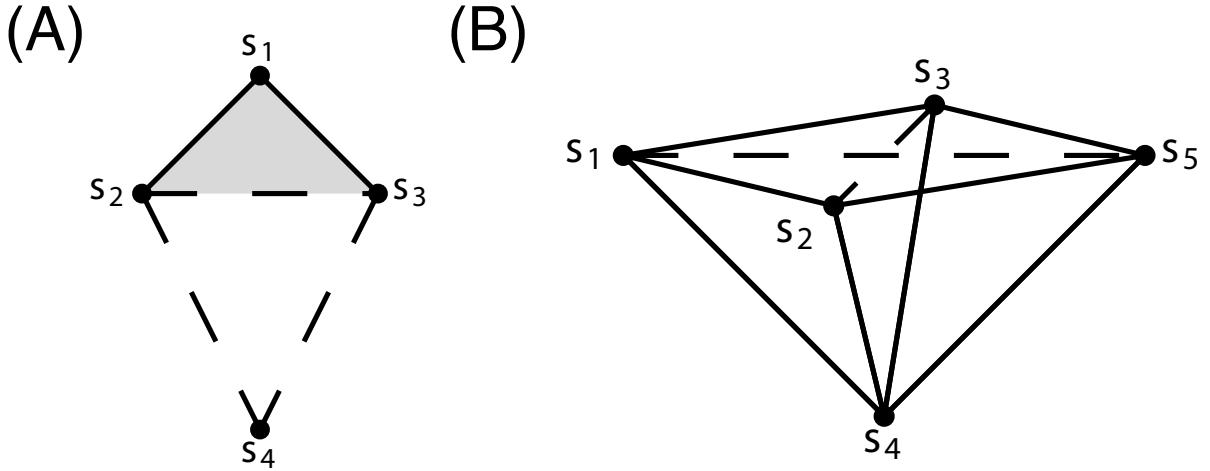


Figure 3.1: Two examples in which the standard filtration fails to identify reticulate evolution. (A) In this example, the ancestral sequences collapse before forming a loop with the recombinant offspring. (B) In this example, multiple recombinations interact to create a degeneracy, and the entire complex collapses immediately. (From Song and Hein [137])

### 3.3 The Median Complex Construction

The median complex is an alternative construction on sequence data aimed at recovering signal of phylogenetic incompatibility using homology. First, we define the median of a set of aligned sequences.

**Definition 1.** For any three aligned sequences  $a$ ,  $b$ , and  $c$ , the *median* sequence  $m(a, b, c)$  is defined such that each position of the median is the majority consensus of the three sequences.

For example, consider the three sequences  $a = 110$ ,  $b = 011$ , and  $c = 101$ . At each site we have the set  $\{1, 1, 0\}$ . The majority consensus for each site is 1, therefore the median sequence is  $m = 111$ . In any further analysis, we augment the original data to include the computed median sequence. Note that as defined here, the median operation is defined only for binary sequences.

Having defined the median operation, we now define the *median closure*. Given an

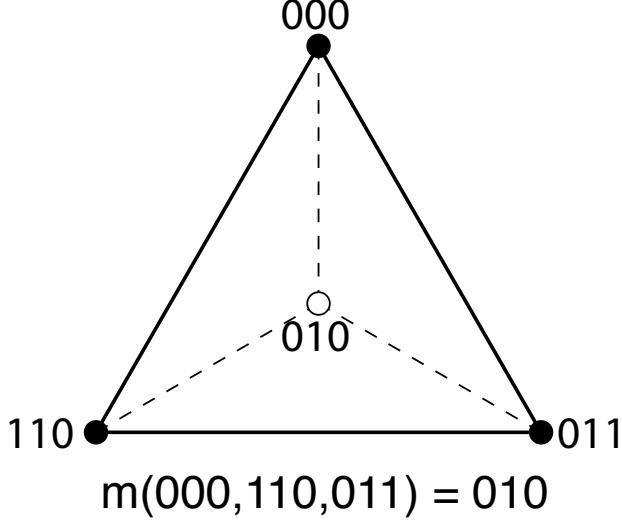


Figure 3.2: The median is defined as the majority allele at each position. The median closure imputes the median into the original vertex set.

alignment  $S$ , the median closure,  $\bar{S}$ , is defined as the vertex set generated from the original set  $S$  that is closed under the median operation,

$$\bar{S} = \{v: v = m(a, b, c) \in \bar{S} \forall a, b, c \in \bar{S}\} \quad (3.1)$$

We can obtain the median closure  $\bar{S}$  by repeatedly applying the median operation to sets of three sequences until no new sequences are added. Effectively, computing the median closure imputes interior nodes into the dataset. We call complexes formed from the original sequences the *leaf complexes*, and call complexes formed from the median closure the *median complexes*. We can then proceed by computing the persistent homology of this median closure. The downside of the median closure operation is that we can no longer identify the loops we measure as reticulate events. The median closure operation can generate multiple loops from a single incompatibility. We now revisit our two examples.

*Example 1.* One median vertex,  $m(s_2, s_3, s_4) = 110$ , as shown in Figure 3.3. This vertex, labeled  $s_r$ , acts as the recombinant offspring of  $s_2$  and  $s_3$ . Persistent homology now detects an  $H_1$  loop in the range  $\epsilon = [1, 2]$  formed between  $s_1$ ,  $s_2$ ,  $s_3$ , and  $s_r$ .  $s_4$  is interpreted the descendant of  $s_r$ .

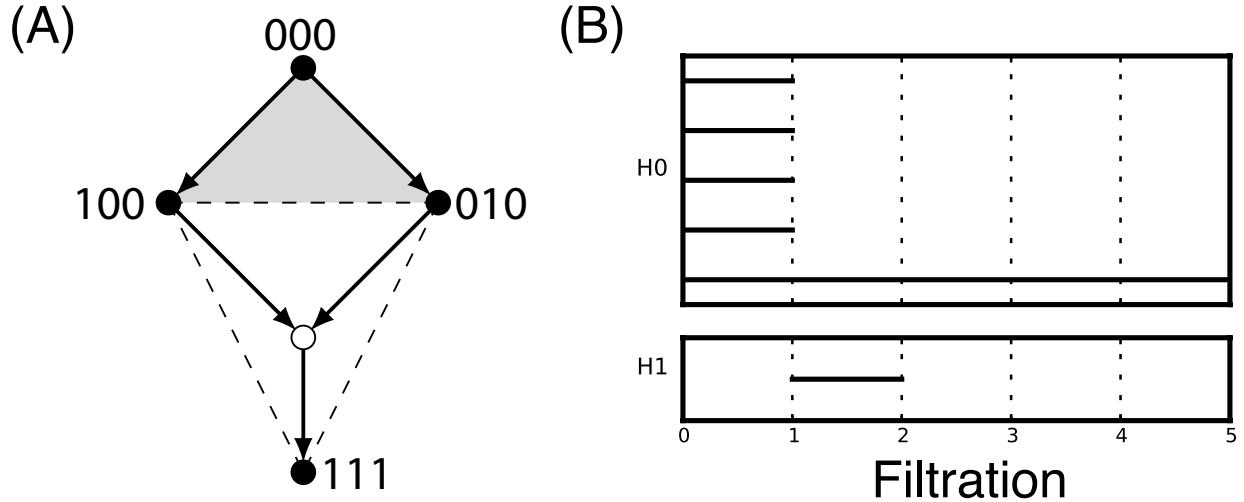


Figure 3.3: One median node (white node), which acts as the recombinant offspring of  $s_2$  and  $s_3$ . One  $H_1$  loop detected in the interval  $[1, 2]$ .

*Example 2.* Four median vertices, as shown in Figure 3.4. Persistent homology now detects four  $H_1$  intervals in the range  $\epsilon = [1, 2)$ . In this case, the median closure now overestimates the minimum number of recombinations required. This example shows a potentially complicating aspect of the median closure in that specific  $H_1$  features are no longer identifiable with specific reticulate events.

Filtrations on Buneman graphs have been defined previously [50], but not using an explicit sequence representation. They have been defined in terms of the split decomposition, which is a deconstruction of the data into sets of possibly-conflicting bipartitions.<sup>1</sup> The filtration defined in Dress, Huber, and Moulton [50] is based on a complicated polytope construction scheme defined directly from the split decomposition. Given that all median graphs are split networks [83], the constructions are identical but the extracted information is not. To the best of our knowledge, quantification of the complexity of these objects has not been measured using homological tools.

---

<sup>1</sup>In a tree, each edge defines a bipartition, or split. A reticulate history will be characterized by incompatible splits.

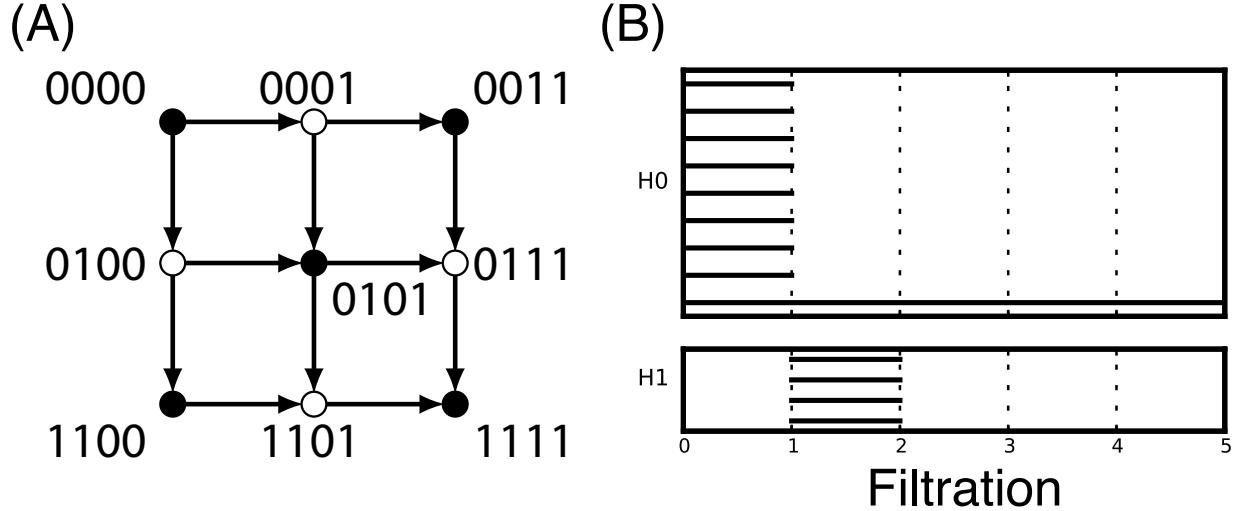


Figure 3.4: Four median vertices (white nodes). Four  $H_1$  loops now detected in the interval  $[1, 2)$ .

### 3.3.1 Inclusion

We have examined the persistent homology of two topological constructions on sequence data: the leaf complex and the median complex. Counting  $\beta_1$  intervals in the leaf complex underestimates reticulate evolution because of incomplete sampling, while counting  $\beta_1$  intervals in the median complex overestimates reticulate evolution. The median complex is in some sense an upper bound on probable recombination histories, and contains within it all possible recombination graphs within it (not strictly true, as there are infinitely many complicated ARGs - but it does contain within it all maximum parsimony trees). We can hypothesize that there exists a true complex, called the *evolutionary complex*, which will accurately reflect the evolutionary relationships in the sequences. Information about the evolutionary complex is not available to us, however we can say that there exists an inclusion between the homotopy types of the three complexes

$$\text{Cl}(\mathcal{LC}) \hookrightarrow \text{Cl}(\mathcal{EC}) \hookrightarrow \text{Cl}(\mathcal{MC}) \quad (3.2)$$

Recovery of an optimal  $\mathcal{EC}$  is the task of many ARG-based methods and is known to be

an NP-hard problem and is not considered here. For example, given an  $\mathcal{EC}$  as computed from some other tool, we might be able to say something useful about the topological complexity.

### 3.3.2 Phylogenetic Examples

Here we consider two standard datasets from the phylogenetics literature. In both examples, the standard filtration yielded no higher homology. We generated the median closure and computed homology on that. Datasets are represented using a triangle-free network construction, which approximates the computed homology.

#### 3.3.2.1 *D. melanogaster* Data

A benchmark dataset in studying recombination is the Kreitman data [96]. The dataset consists of eleven sequences (nine unique) of the Adh locus from *Drosophila melanogaster* collected from various geographic locations, with 43 segregating sites. The Hudson-Kreitman test yields 6 reticulate events. Computing the median closure expands the dataset to 46 vertices. Here we have non-trivial homology: 32  $H_1$  loops and 3  $H_3$  loops. In the visualized network, the complex reticulations ( $H_3$ ) are localized to the bottom-most samples. The  $H_1$  reticulations, on the other hand, are not very localized and persist across geographic regions. The barcode plot is shown in Figure 3.5.

#### 3.3.2.2 *Ranunculus* Data

Natural hybridization occurs frequently in plants. Here we examine reticulation in the maturase K (matK) protein in nine species from genus *Ranunculus*. This data is originally from [79]. From nine initial species, the median closure has 32 vertices. Persistent homology is computed and the barcode diagram shown in Figure 3.6. Looking at  $H_0$ , we identify two clusters of species. Further, we identify 17  $H_1$  loops and 3  $H_3$  loops. Comparing with the *D. melanogaster* data, reticulation at this locus is both smaller in scale (shorter bars at small

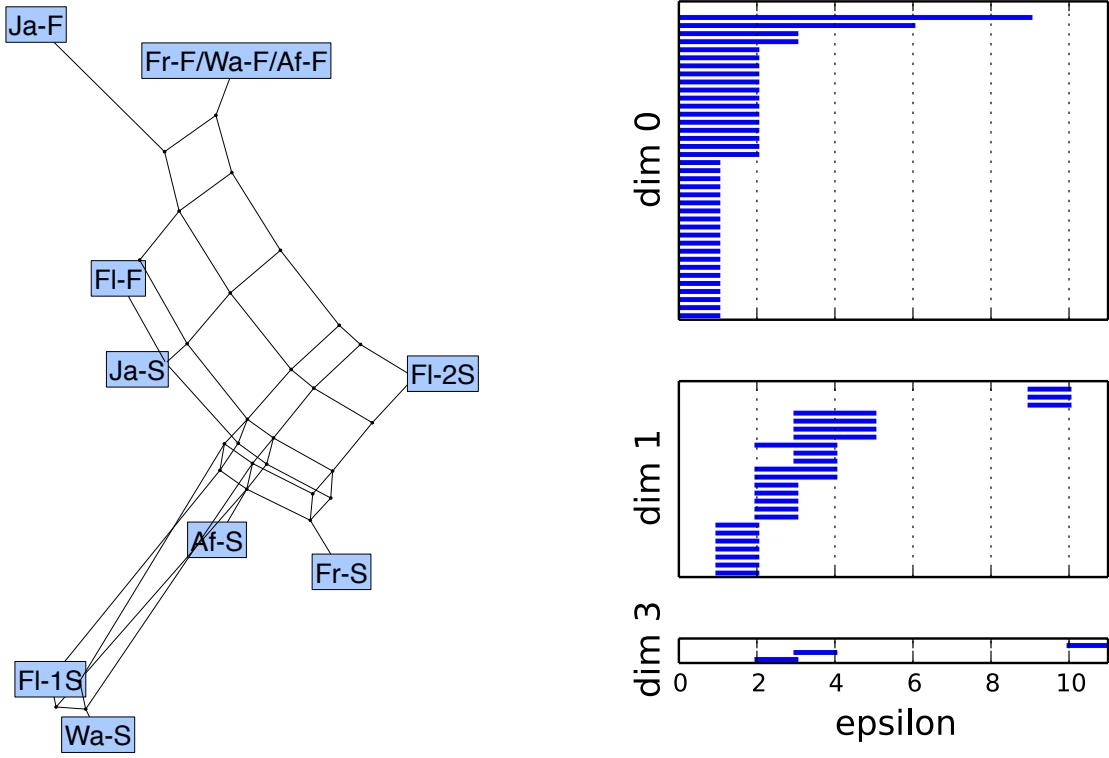


Figure 3.5: Recombination in *D. melanogaster*. Persistent homology identifies several complex reticulations in the population.

filtration values) and less frequent (fewer total bars). Additionally, the complex reticulations are localized within each  $H_0$  cluster.

## 3.4 Čech Complex Construction as an Optimization Problem

The Čech complex is defined on a set of points  $S$  as

$$\check{\text{C}}\text{ech}(r) = \left\{ \sigma \subseteq S \mid \bigcap_{x \in \sigma} B_x(r) \neq \emptyset \right\}, \quad (3.3)$$

where  $B_x(r)$  is the ball of radius  $r$  centered at vertex  $x$ . By the nerve lemma, the homotopy type of the Čech covering is guaranteed to be identical to that of the original topological

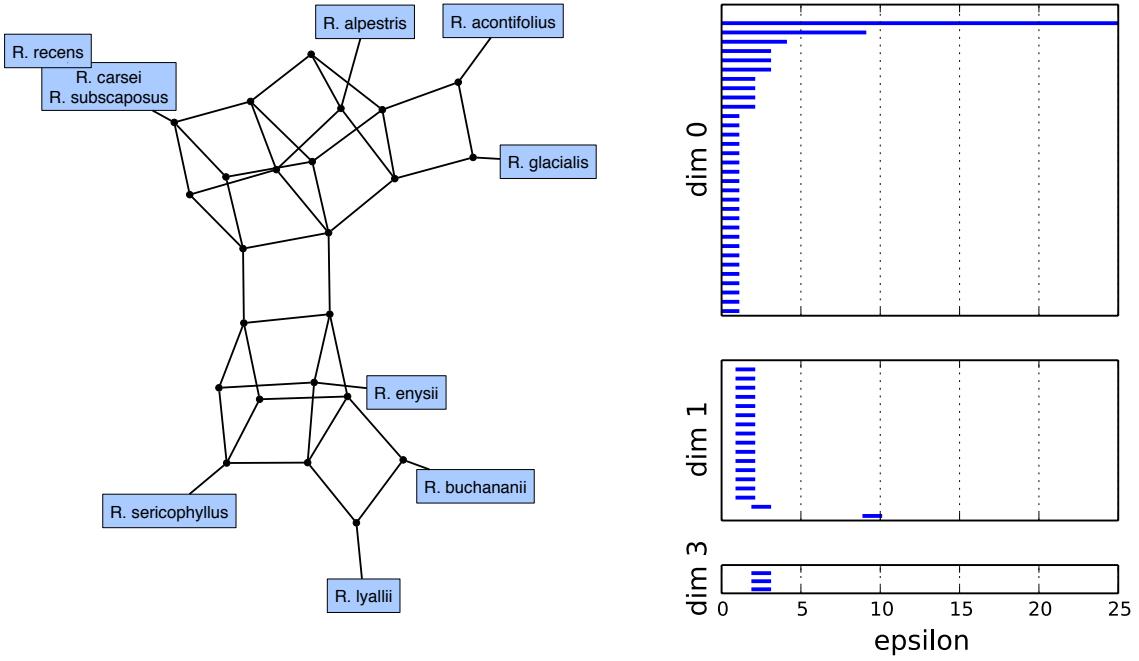


Figure 3.6: Species hybridization in genus *Ranunculus*. Persistent homology identifies two populations separated by complex reticulations.

space [17].

Computing the Čech complex is often an expensive operation, such that in practice the Vietoris-Rips complex is used. Unlike the Vietoris-Rips complex, which is entirely defined by the 1-skeleton, the Čech complex requires one to check each simplex  $\sigma$  up to some maximum dimension  $D$ . The Čech complex therefore requires one to know the ambient space the data is embedded in, unlike a Rips complex which can be built directly from distance data. Binary sequence data of length  $d$  explicitly sits on the discrete lattice of  $\{0, 1\}^d$  with an  $L_1$  norm. In this case, it is not immediately obvious how to define when three sequences should form a simplex. One Therefore, we expand the ambient space to  $\mathbb{R}^d$  with an  $L_1$  metric. This choice of metric is motivated by two reasons. First, the  $L_1$  norm maintains the Hamming distance between sampled points. Second, the  $L_1$  norm keeps the primary theorem intact,

that is tree like data generates trivial homology.<sup>2</sup>

The problem of deciding if a particular simplex  $\sigma$  belongs in the Čech complex at radius  $r$  is the same as checking if a ball of radius  $r$  can be placed such that each point  $x$  in  $\sigma$  is contained within the ball. In  $\mathbb{R}^d$  with an  $L_2$  metric there exists an efficient randomized algorithm for computing this radius known as the *miniball algorithm*.<sup>[68]</sup> However, the efficiency of the miniball algorithm relies on the strict convexity of the  $L_2$  metric and therefore is not applicable to a space with an  $L_1$  metric. Instead, we pose the miniball problem in  $L_1$  as a generic convex optimization problem, and use standard library solver. That is, we define a  $d+1$  dimensional optimization problem where  $x$  is the miniball center and  $R$  is the miniball radius.

The problem is stated as

$$\begin{aligned} & \text{minimize} && R \\ & \text{subject to} && \forall p \in P : \|x - p\|_1 \leq R \\ & && x \in \mathbb{R}^d \end{aligned}$$

We implement the optimization problem in `cvxpy`.

### 3.4.1 Molecular Hypothesis

Gromov proved that a median graph is the 1-skeleton of a CAT(0) cubical complex [74]. The homology of a cubical complex can be efficiently computed using the methods of Kaczynski, Mischaikow, and Mrozek [89] through a slightly different construction. We define a cubical flag complex and build a filtration dimension by dimension (to expand on this point...) The barcode diagram will then have the natural interpretation of being composed of sets of hypercubes of varying dimension. If we consider each bar of dimension  $n$  in the barcode diagram in turn, we can determine the incompatible sites that it represents. Dimension 1 bars (2-cubes) will have one pair of incompatible sites with four haplotypes. Dimension 2

---

<sup>2</sup>This notion has a natural extension to multiallelic sites which is not detailed here.

Table 3.1: Čech Homology of Hypercube

$d =$	1	2	3	4	5	6
$H_0$	2	4	8	16	32	64
$H_1$	0	1	5	17	49	129
$H_2$	0	0	1	7	31	111
$H_3$	0	0	0	1	9	49
$H_4$	0	0	0	0	1	11
$H_5$	0	0	0	0	0	1
$H_6$	0	0	0	0	0	0

bars (3-cubes) will have three pairs of incompatible sites with eight haplotypes. In general,  $n$  bars will represent  $n + 1$ -cubes in which all  $2^{(n+1)}$  haplotypes are present in the vertices of the generating cycle.

From the barcode diagram it will not in general be possible to decompose our construction into the primitive building blocks of hypercubes. This is because the hypercubes of dimension ( $n > 2$ ) will in general not be independent, but can interact by sharing lower dimensional faces. Nonetheless, to aid in decomposing the barcode diagram, we constructed the following table, which contains the homology ranks (betti numbers) for powers of the hypercube graph, computed using the Čech complex. Incidentally, it was understanding the structure of numbers in a table very much like Table 3.1 which led us to find a method of computing Čech homology instead of Rips homology.

## 3.5 Conclusions

Persistent homology can capture and quantify complex patterns of reticulation in genomic data. The standard Vietoris-Rips filtration is susceptible to reduced sensitivity due to incomplete sampling or interactions between reticulations. Constructing the median closure of the original sequence set increases the topological signal of reticulation. Future work will focus on efficient implementations of constructing this closure. We also introduced a Čech complex construction on genomic data. The construction treats filling higher-dimensional

simplices as an optimization problem, which is solved using the miniball algorithm. An interesting additional observation is that the number of recombinations required to explain the fully saturated hypercube is exactly equal to the alternating sum of the homology ranks.



# Chapter 4

## Parametric Inference using Persistence Diagrams

*“I predict a new subject of statistical topology. Rather than count the number of holes, Betti numbers, etc., one will be more interested in the distribution of such objects on noncompact manifolds as one goes out to infinity”*

*Isadore Singer*

### 4.1 Introduction

Recent work in topological data analysis has concentrated on developing the statistical foundations for data analysis using the persistent homology framework (see the discussion in Section 2.2.2.2 and references [62, 16, 34]). The focus of this work has primarily been estimating the topology of an object from a finite, noisy sample. Doing so requires statistical methods to distinguish topological signal from noise.

Here we consider a different scenario. Many simple stochastic models generate complex data that cannot be readily visualized as a manifold or summarized by a small number of topological features. The persistence diagrams generated from such models will be unique in two ways: (1) the complexity of the diagram (i.e. number of topological features) will grow

with the number of sampled points, and (2) each instantiation of the model will generate a unique set of topological features. Nevertheless, the collection of measured topological features may exhibit additional structure, providing useful information about the underlying data generating process. While the persistence diagram is itself a summary of the topological information contained in a sampled point cloud, to perform inference further summarization may be appropriate, e.g. by considering distributions of properties defined on the diagram. In other words, we are less interested in learning the topology of a particular sample, but rather in understanding the expected topological signal of different model parameters. We show an example in Figure 4.1. Here we have three identical simulations of a stochastic coalescent model, commonly used in population genetics. Each simulation is generated with the same number of points and the same parameters. Because the model is random, the resulting topology will also be random, making it impossible to match individual topological features between diagrams. We would like a way to characterize these models using topology.

In this chapter, we show that summary statistics computed on the persistence diagram can be used for likelihood-based parametric inference. We use genomic sequence data as a case study, examining the topological behavior of the coalescent process with recombination, a widely used stochastic model of biological evolution. We find that the process generates nontrivial topology in a way that depends sensitively on parameters in the model. The idea is presented as a proof of concept, in order to motivate the identification additional models with regular topological structure that may amenable to this type of inference.

There has been related work on the topology of random models, mostly in the non-persistent case. The statistical properties of random simplicial complexes, including distributions over their Betti numbers, has been studied in [90, 91]. The persistent homology of Gaussian random fields and other probabilistic structures has been studied in [1]. Functions defined on the persistence diagram were used to compute a fractal dimension for various polymer physics models in [109].

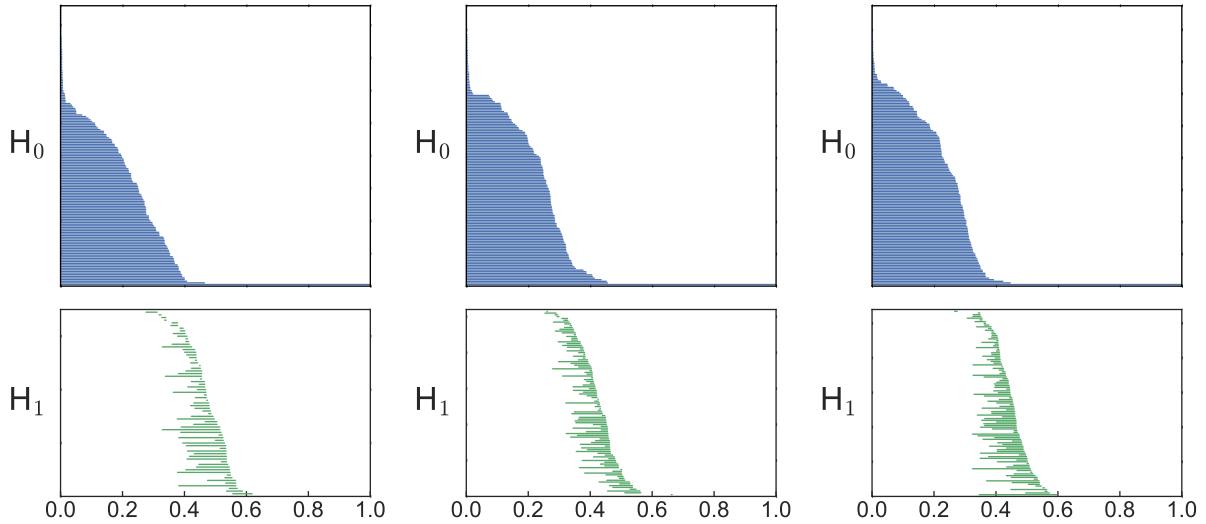


Figure 4.1: Barcode diagrams for three simulations of a stochastic coalescent model. Each simulation is generated using the same parameters, however as random models the persistent homology will differ for each run. Simulations were performed in `ms` using the command `ms 200 1 -t 500 -r 144 10000`.

## 4.2 Warmup: Gaussian Random Fields

Here we show that parametric inference of Gaussian Random Fields can be performed from the barcode diagram. Make reference to [1]. Connections with problems in cosmology.

**TODO:** expand

## 4.3 The Coalescent Process

The coalescent process is a stochastic model that generates the genealogy of individuals sampled from an evolving population [148]. The genealogy is then used to simulate the genetic sequences of the sample. This model is essential to many methods commonly used in population genetics. Starting with a present-day sample of  $n$  individuals, each individual's lineage is traced backward in time, towards a mutual common ancestor. Two separate lineages collapse via a coalescence event, representing the sharing of an ancestor by the two lineages. The stochastic process ends when all lineages of all sampled individuals collapse

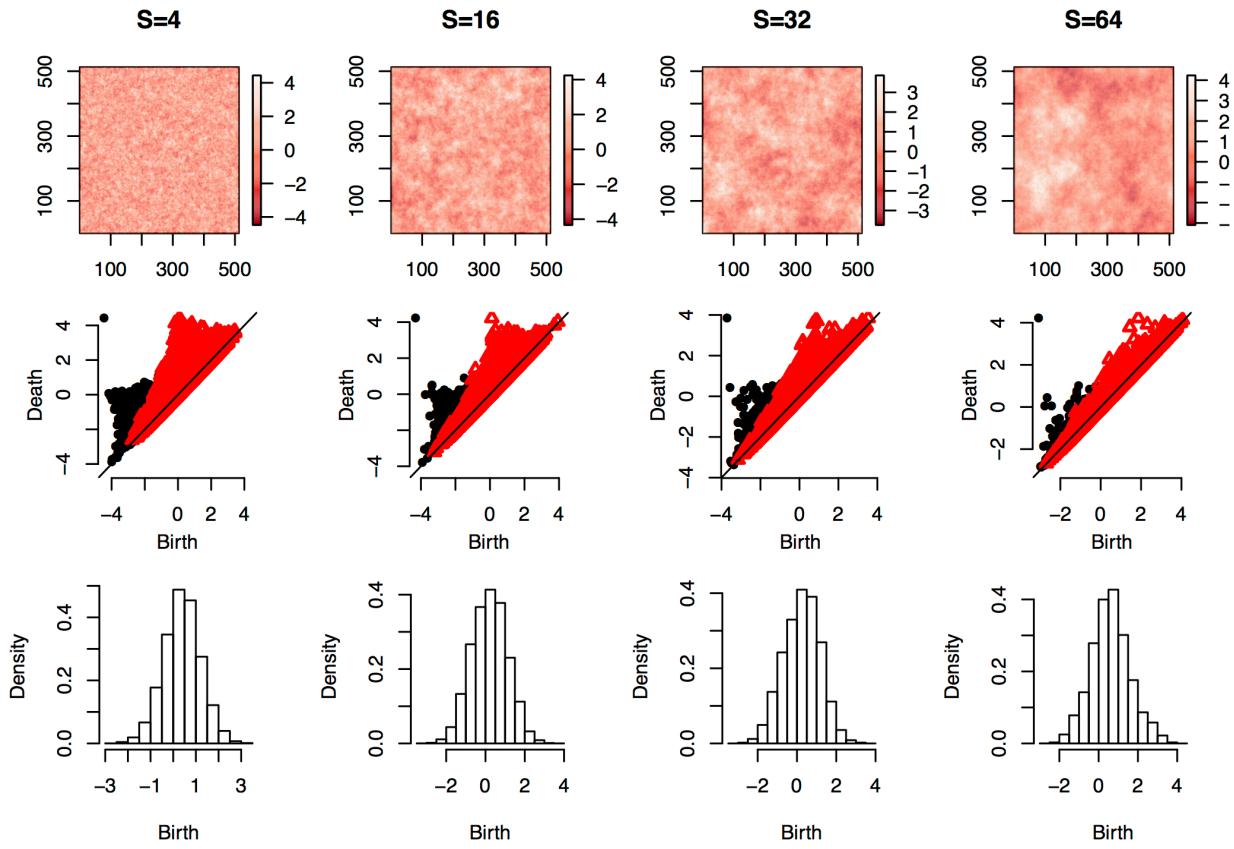


Figure 4.2: Gaussian Random Fields with Exponential Covariance Structure.

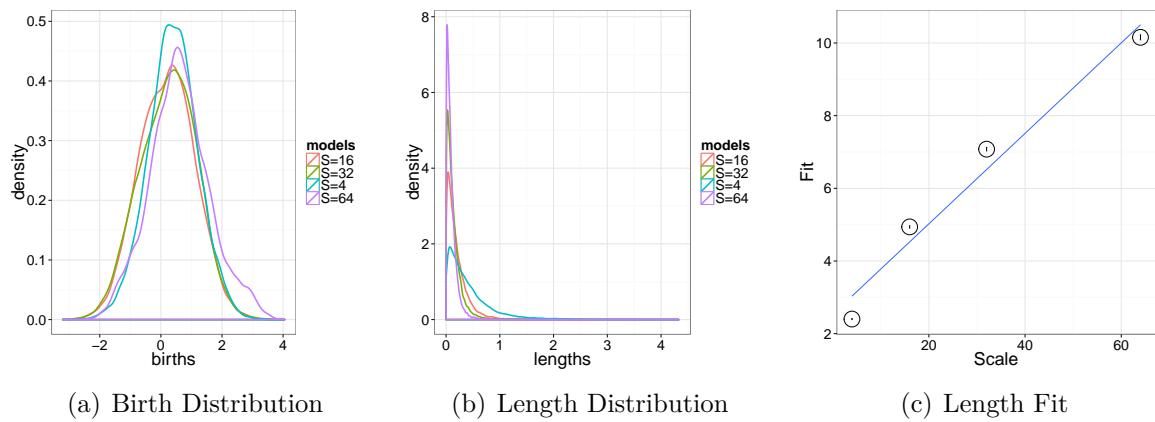


Figure 4.3: Gaussian Random Field Summary Statistics

into a single common ancestor. In this process, if the total (diploid) population size  $N$  is sufficiently large, then the expected time before a coalescence event, in units of  $2N$  generations, is approximately exponentially distributed:

$$P(T_k = t) \approx \binom{k}{2} e^{-(\frac{k}{2})t}, \quad (4.1)$$

where  $T_k$  is the time that it takes for  $k$  individual lineages to collapse into  $k - 1$  lineages.

After generating a genealogy, the genetic sequences of the sample can be simulated by placing mutations on the individual branches of the lineage. The number of mutations on each branch is Poisson-distributed with mean  $\theta t/2$ , where  $t$  is the branch length and  $\theta$  is the population-scaled mutation rate. In this model, the average *genetic distance* between any two sampled individuals, defined by the number of mutations separating them, is  $\theta$ .

The coalescent with recombination is an extension of this model that allows different genetic loci to have different genealogies. Looking backward in time, recombination is modeled as a splitting event, occurring at a rate determined by population-scaled recombination rate  $\rho$ , such that an individual has a different ancestor at different loci. Evolutionary histories are no longer represented by a tree, but rather by an *ancestral recombination graph*. Recombination is the component of the model generating nontrivial topology by introducing deviations from a contractible tree structure, and is the component which we would like to quantify. Coalescent simulations were performed using `ms` [81].

## 4.4 Statistical Model

The persistence diagram from a typical coalescent simulation is shown in Figure 4.1. Examining the diagram, it would be difficult to classify the observed features into signal and noise. Instead, we use the information in the diagram to construct a statistical model in order to infer the parameters,  $\theta$  and  $\rho$ , which generated the data. Note that we consider inference using only  $H_1$  invariants, but the ideas easily generalize to higher dimensions. We consider the following properties of the persistence diagram: the total number of features,  $K$ ; the set

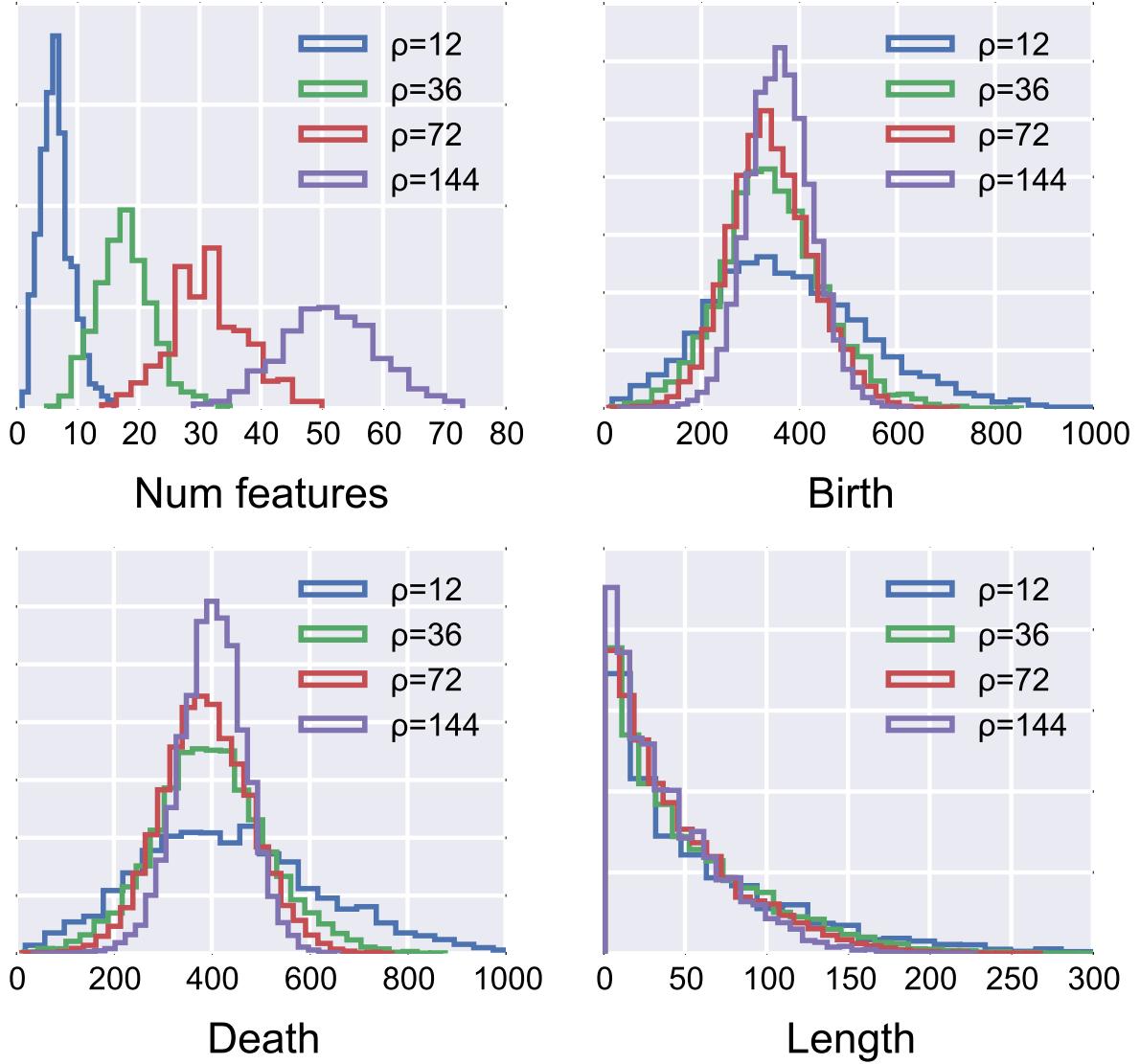


Figure 4.4: Distributions of statistics defined on the  $H_1$  persistence diagram for different model parameters. Top left: Number of features. Top right: Birth time distribution. Bottom left: Death time distribution. Bottom right: Feature length distribution. Data generated from 1000 coalescent simulations with  $n = 100$ ,  $\theta = 500$ , and variable  $\rho$ .

of birth times,  $(b_1, \dots, b_K)$ ; the set of death times,  $(d_1, \dots, d_K)$ ; and the set of persistence lengths,  $(l_1, \dots, l_K)$ . In Figure 4.4 we show the distributions of these properties for four values of  $\rho$ , keeping fixed  $n = 100$  and  $\theta = 500$ . Several observations are immediately apparent. First, the topological signal is remarkably stable. Second, higher  $\rho$  increases the number of features, consistent with the intuition that recombination generates nontrivial topology in the model. Third, the mean values of the birth and death time distributions are only weakly dependent on  $\rho$  and are slightly smaller than  $\theta$ , suggesting that  $\theta$  defines a natural scale in the topological space. However, higher  $\rho$  tightens the variance of the distributions. Finally, persistence lengths are independent of  $\rho$ .

Examining Figure 4.4, we can postulate:  $K \sim \text{Pois}(\zeta)$ ,  $b_k \sim \text{Gamma}(\alpha, \xi)$ , and  $l_k \sim \text{Exp}(\eta)$ . Death time is given by  $d_k = b_k + l_k$ , which is incomplete Gamma distributed. The parameters of each distribution are assumed to be an *a priori* unknown function of the model parameters,  $\theta$  and  $\rho$ , and the sample size,  $n$ . Keeping  $n$  fixed, and assuming each element in the diagram is independent, we can define the full likelihood as

$$p(D | \theta, \rho) = p(K | \theta, \rho) \prod_{k=1}^K p(b_k | \theta, \rho) p(l_k | \theta, \rho). \quad (4.2)$$

Simulations over a range of parameter values suggest the following functional forms for the parameters of each distribution. The number of features is Poisson distributed with expected value

$$\zeta = a_0 \log \left( 1 + \frac{\rho}{a_1 + a_2 \rho} \right) \quad (4.3)$$

Birth times are Gamma distributed with shape parameter

$$\alpha = b_0 \rho + b_1 \quad (4.4)$$

and scale parameter

$$\xi = \frac{1}{\alpha} (c_0 \exp(-c_1 \rho) + c_2). \quad (4.5)$$

These expressions appears to hold well in the regime  $\rho < \theta$ , but break down for large  $\rho$ . The length distribution is exponentially distributed with shape parameter proportional

to mutation rate,  $\eta = \alpha\theta$ . The coefficients in each of these functions are calibrated using simulations, and could be improved with further analysis. This model has a simple structure and standard maximum likelihood approaches can be used to find optimal values of  $\theta$  and  $\rho$ .

## 4.5 Coalescent Simulations

We simulated a coalescent process with sample size  $n = 100$  and  $l = 10,000$  loci. The mutation rate,  $\theta$ , was varied across  $\theta = \{50, 500, 5000\}$ . The recombination rate,  $\rho$ , was varied across  $\rho = \{4, 12, 36, 72\}$ . The output of the process is a set of binary sequences of variable length (length is dependent on  $\theta$ ). The Hamming metric is used to construct a pairwise distance matrix between sequences. We computed persistent homology and used the model described in Section 4.4 to estimate  $\theta$  and  $\rho$ . Results are shown in Figure 4.5, where we plot estimates and 95% confidence interval from 500 simulations. We observe an improved  $\rho$  estimate at higher mutation rate. This is expected, as increasing  $\theta$  is essentially increasing sampling on branches in the genealogy. We also observe tighter confidence intervals at higher recombination rates, consistent with the behavior seen in Figure 4.4.

As a final point of discussion, we make a comparison with existing methods of estimating the recombination rate in coalescent models. One of the earliest methods was developed by Hudson in [80] and relies on a modeling the distribution of pairwise distances in sequence data. Our estimator, by contrast, is built on modeling the distribution of topological features measured in the data. In Figure 4.6 we show what these distributions look like for a few values of  $\rho$ . The distribution of pairwise distances follows an interesting pattern as we increase the  $\rho$ . At  $\rho = 0$ , the distances follow an exponential distribution, and as  $\rho$  is increased the distances begin to follow a normal distribution with an increasingly tight variance. That is to say, as  $\rho$  is increased, all pairs of sequences begin to be normally distributed around some mean value (that is determined by the mutation rate). Estimating  $\rho$  from this data

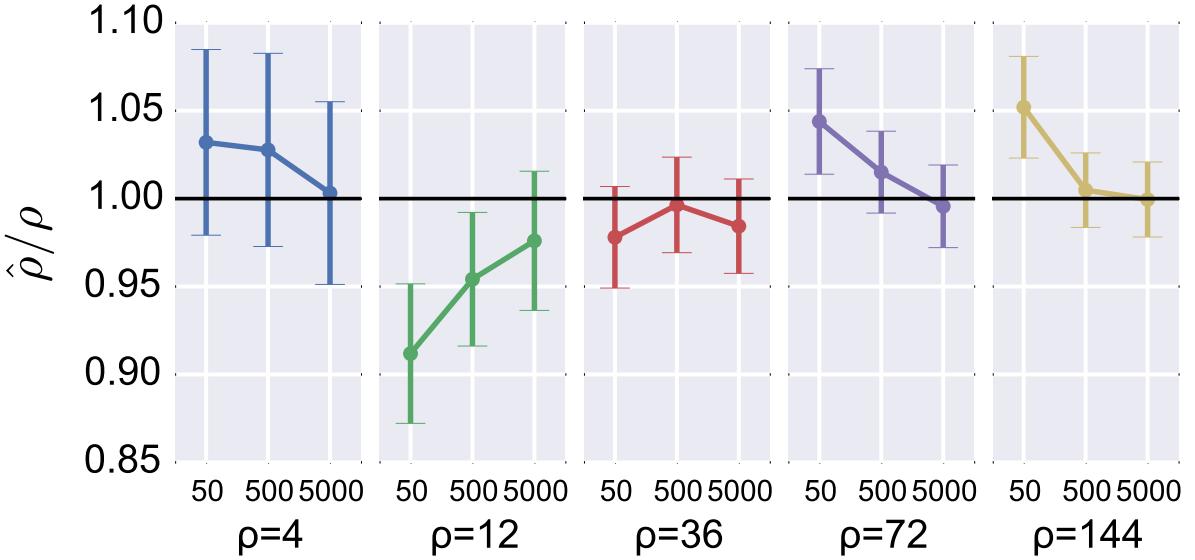


Figure 4.5: Inference of recombination rate  $\rho$  using topological information. The recombination rate  $\rho$  is estimated for five values  $\{4, 12, 36, 72, 144\}$  at three different mutation rates  $\{50, 500, 5000\}$ . Mean estimate over 500 simulations and 95% confidence interval is shown.

requires teasing out the contribution of recombination from the contribution of coalescence. In contrast, our topological estimator is in some sense a more pure signal of recombination. By the fundamental theorem, there will be no  $H_1$  homology when  $\rho = 0$ . Any signal that is generated at  $H_1$  and higher is due strictly to reticulate processes.

## 4.6 Conclusions

In machine learning, the task is often to infer parameters of a model from observations. In this chapter we have presented a proof of concept for statistical inference based on topological information computed using persistent homology. Unlike previous work, which considered estimating homology of a partially observed object, we were interested in a model which generates a complex, but stable, topological signal. Three conditions were required for the success of this approach: First, a well-defined statistical model. Second, an intuition that the observed topological structure is directly correlated with the parameters of interest in the

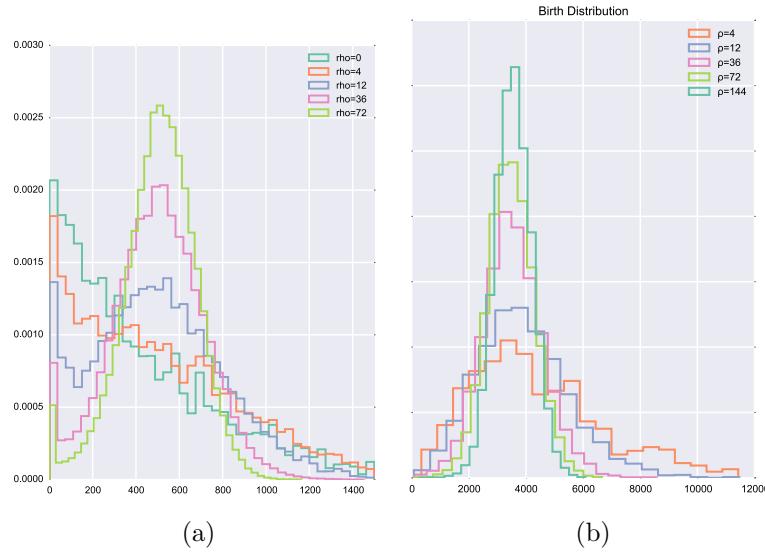


Figure 4.6: Traditional estimators of  $\rho$  are based on modeling the distribution of pairwise distances in a dataset, as shown on left. This distribution is a mixture of both coalescent and recombination processes. In contrast, the distribution of birth times is driven solely by recombination.

model. Third, sufficient topological signal to reliably estimate statistics on the persistence diagram. It is an open question to identify classes of models for which these conditions will hold.

# **Part II**

# **Applications**



# Chapter 5

## Phage Mosaicism

### 5.1 Introduction

Phages are microbial viruses which can infect bacteria, archaea, or single-celled eukaryotes. By some measures, they are the most abundant and diverse class of organism on the planet. It is estimated that there are  $10^{31}$  extant bacteriophages [132].<sup>1</sup> The phage population completely turns over every few days – an estimated infection rate of  $10^{23}$  per second [140].

Phages play an essential role in natural ecosystems by regulating bacterial populations. Steps have been taken towards harnessing this ability for productive use – the FDA has approved several bacteriophage products designed to kill harmful bacteria in dairy and meat products [21]. Also promising are potential phage therapies for treating pathogenic bacterial infections, although research in this direction is controversial [93].

Phages are classified based on lifestyle: virulent phages have a lytic life cycle and will infect a host, multiply, and exit the cell via lysis, killing the host organism; temperate phages have a lysogenic life cycle and can remain within the host in a latent state, without disrupting host cellular function. Phages can have a nucleic acid composition that is either double-

---

<sup>1</sup>The estimate can be arrived at two independent ways: by assuming a total bacterial population size of  $10^{30}$ , and approximately ten phages per bacteria; or by the observation of a phage density of  $10^6$  to  $10^7$  per mL of seawater.

Table 5.1: Phage families defined by the ICTV

Order	Family	Morphology	Nucleic acid
<i>Caudovirales</i>	<i>Myoviridae</i>	Nonenveloped, contractile tail	linear dsDNA
	<i>Siphoviridae</i>	Nonenveloped, noncontractile tail (long)	linear dsDNA
	<i>Podoviridae</i>	Nonenveloped, noncontractile tail (short)	linear dsDNA
<i>Ligamenvirales</i>	<i>Lipothrixviridae</i>	Enveloped, rod-shaped	linear dsDNA
	<i>Rudiviridae</i>	Nonenveloped, rod-shaped	linear dsDNA
Unassigned	<i>Ampullaviridae</i>	Enveloped, bottle-shaped	linear dsDNA
	<i>Bicaudaviridae</i>	Nonenveloped, lemon-shaped	circular dsDNA
	<i>Clavaviridae</i>	Nonenveloped, rod-shaped	circular dsDNA
	<i>Corticoviridae</i>	Nonenveloped, isometric	circular dsDNA
	<i>Cystoviridae</i>	Enveloped, spherical	segmented dsRNA
	<i>Fuselloviridae</i>	Nonenveloped, lemon-shaped	circular dsDNA
	<i>Globuloviridae</i>	Enveloped, isometric	linear dsDNA
	<i>Guttaviridae</i>	Nonenveloped, ovoid	circular dsDNA
	<i>Inoviridae</i>	Nonenveloped, filamentous	circular ssDNA
	<i>Leviviridae</i>	Nonenveloped, isometric	linear ssRNA
	<i>Microviridae</i>	Nonenveloped, isometric	circular ssDNA
	<i>Plasmaviridae</i>	Enveloped, pleomorph	circular dsDNA
	<i>Tectiviridae</i>	Nonenveloped, isometric	linear dsDNA

stranded DNA (dsDNA), single-stranded DNA (ssDNA), double-stranded RNA (dsRNA), or single-stranded RNA (ssRNA). Of these, dsDNA is by far the most common. The typical phage genome length is on the order of  $10^5$  bases, but can range from  $10^3$  to  $10^6$  bases.

Because there is no conserved gene across all phage populations, there is no accepted way of constructing a molecular phage taxonomy. The current bacteriophage taxonomy is compiled by the International Committee on Taxonomy of Viruses (ICTV) and is based on virus morphology, host range, lifestyle, and nucleic acid composition [86]. Table 5.1 presents an overview of phage families as defined by the ICTV. There are two assigned orders and eighteen recognized families. Fourteen families have dsDNA, two families have ssDNA, and two families have an RNA genome.

Phages have been shown to be subject to high rates of reticulate genomic exchange [152]. The phage genome was believed to be mosaic, composed of distinct modules that can be freely exchanged within a population. Increased genomic data has confirmed this mosaic structure and raised questions about the applicability and interpretation of the ICTV taxonomy. Based solely on morphology and host, the ICTV taxonomy has been shown to be inconsistent with the genomic data, as the following example from Lawrence *et al.* shows [101]. In Figure 5.1

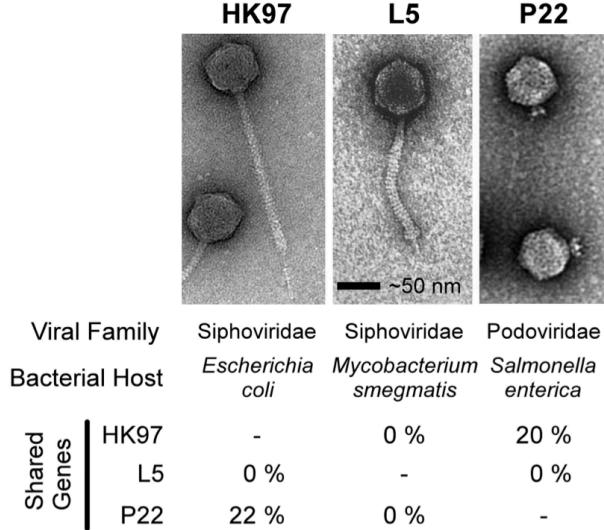


Figure 5.1: Inconsistency of morphological classifications in bacteriophage. HK97 and L5 are classified in the Siphoviridae family of long tail non-contractile phages, despite sharing no gene content. P22, a short-tail phage in the Podoviridae family, while morphologically dissimilar, shares 20% gene content with HK97. Figure adapted from [101].

we show three different bacteriophage species: Enterobacteria phage HK97, Mycobacterium phage L5, and Enterobacteria phage P22. HK97 is a Siphoviridae infecting *E. coli*. L5 is a Siphoviridae infecting *M. smegmatis*. P22 is a Podoviridae infecting *S. enterica*. HK97 and L5 belong to the Siphoviridae family comprised of long tail noncontractile phages. P22 belongs to the Podoviridae family comprised of short tail phages. Visually, it appears that HK97 and L5 should indeed be classified as distinct from P22. However, genomic analysis reveals that HK97 and L5 share no gene content, and, despite appearances to the contrary, HK97 and P22 share 20% gene content. This example demonstrates that morphology and host range alone are not sufficient in representing phage relationships.

Alternative representations of phage relationships have been proposed based on whole genome analysis. For example, Rohwer and Edwards constructed a phage phylogenetic tree using differences in phage proteomes [131]. Proux *et al.* proposed a phylogenetic representation based on comparative analysis of head and tail sequences [127]. However, these models still make the assumption of tree-like relationships, which will not be appropriate for representing highly mosaic molecular relationships.

In this chapter, we use approaches from topological data analysis to identify, measure, and represent reticulate evolution in a population of phage sequences. This work is primarily based on data collected by Lima-Mendez *et al.* [104]. First, we use persistent homology to characterize reticulation in phage genomes. We find  $H_0$  is largely inconsistent with existing phage taxonomies, and interpret  $H_1$  as evidence for reticulate genetic exchange due to shared ecology and host range. Second, we visualize phage molecular relationships using Mapper, identifying clusters of phages with common gene content and host range. Representative protein families for each phage cluster are identified. The Mapper network suggests an alternate way of representing phage molecular relationships.

## 5.2 Data

We use data initially collected and analyzed in [104]. The initial data set consists of a collection of 306 sequenced bacteriophage genomes. We show summary information about the data in Figure 5.2. Of the 306 genomes, 246 consist of dsDNA, 36 ssDNA, 12 dsRNA, and 8 ssRNA. Four have unclassified nucleic acid material. With respect to lifestyle, 146 are temperate and 72 are virulent. Actinoplanes phage phiAsp2 is the single pseudotemperate phage, which means it largely maintains a temperate lifestyle but can occasionally enter a virulent state. For 87 phages the lifestyle is unknown. Taxonomically, the vast majority belong to order Caudovirales (221), which comprises Siphoviridae (117), Myoviridae (47), and Podoviridae (54). Order Ligamenvirales (4) comprises Lipothrixviriae (2) and Ravidviridae (2). Unassigned families include Inoviridae (22), Cystoviridae (12), Gokushoviridae (8), and Microviridae (6).

Each of the 306 bacteriophage genomes has been sequenced and annotated.<sup>2</sup> This step resulted in 19,537 unique bacteriophage phage genes. In the original study [104], these

---

<sup>2</sup>The annotation step assigns genes to subsequences of the genome. For well-characterized species this is facilitated by a reference genome. For less well-characterized species this can require the use of heuristic gene-finding algorithms.

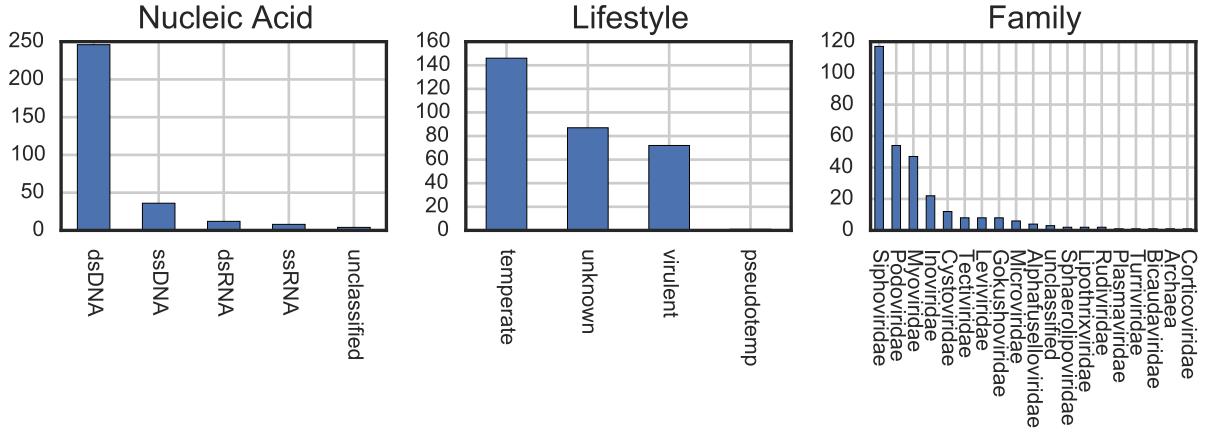


Figure 5.2: Summary annotations of phage data used in this analysis. 306 bacteriophage genomes were included, as originally collected in [104]. Here we show various annotations for the phage, including nucleic acid type, lifestyle, and taxonomic family (as defined by the ICTV). For some phage strains this data is unknown.

genes were then clustered into 8,576 protein families using BlastP, which analyzes pairwise similarity of proteins [3]. Protein families share homology, which implies some degree of shared evolutionary ancestry. Phages can then be represented as phyletic profiles in a protein family-space, indicating the presence or absence of a particular protein family. In this case, the phyletic matrix  $P$  is a  $306 \times 8576$  binary matrix.

## 5.3 Measuring Phage Mosaicism with Persistent Homology

We apply persistent homology to the phyletic profiles in order to quantify reticulation in the bacteriophage data. Because we have transformed from sequence space into phyletic profiles, we do not invoke a specific evolutionary model. However, the fundamental theorem that non-trivial homology implies reticulation still holds. First, we construct an appropriate metric space. Following [104], we use a hypergeometric model as follows. For two phages  $A$  and  $B$ , let  $a$  be the number of protein families in phage  $A$ ,  $b$  be the number of protein

families in phage  $B$ , and  $c$  be the number of protein families in common. Let  $n$  be the total number of protein families. Then we can compute the p-value that the number of shared protein families  $c$  is significant as

$$P_{AB} = \sum_{i=c}^{\min(a,b)} \frac{\binom{a}{i} \binom{n-a}{b-i}}{\binom{n}{b}}. \quad (5.1)$$

To convert the p-values into a distance we take the log transform with small added noise,

$$d_{AB} = \log_{10}(P_{AB} + 10^{-10}) + 10. \quad (5.2)$$

This yields a distance matrix  $D$  with distances scaled between 0 and 10. While this space does not explicitly reflect evolutionary divergence at a molecular level, it may be realistic at the protein level at which more complex types of genome evolution will have occurred.

We now compute the persistent homology of  $D$ . The barcode diagram is shown in Figure 5.3. The  $H_0$  information represents hierarchical clustering and can be identically represented as a dendrogram. We show the dendrogram, restricting only to strains of order Caudovirales, in Figure 5.4. The strains are labeled by their taxonomic family: red for Myoviridae, blue for Siphoviridae, and green for Podoviridae. We can immediately see that the assigned taxonomic families are not consistent with the clustering based on protein information. However, there does appear to be some structure in which the taxonomic label is consistent within clusters of strains. Returning to the barcode diagram, we see substantial nontrivial homology in  $H_1$  across all scales. This confirms the presence of mosaic exchange expected in phage genomes.

Focusing on order Caudovirales, for which the most data was present. We separately computed persistent homology for each of the three families. The barcode diagrams are shown in Figure 5.5. Computing the TOP score for each family, we have Myoviridae = 0.58, Siphoviridae = 1.14, and Podoviridae = 0.56.

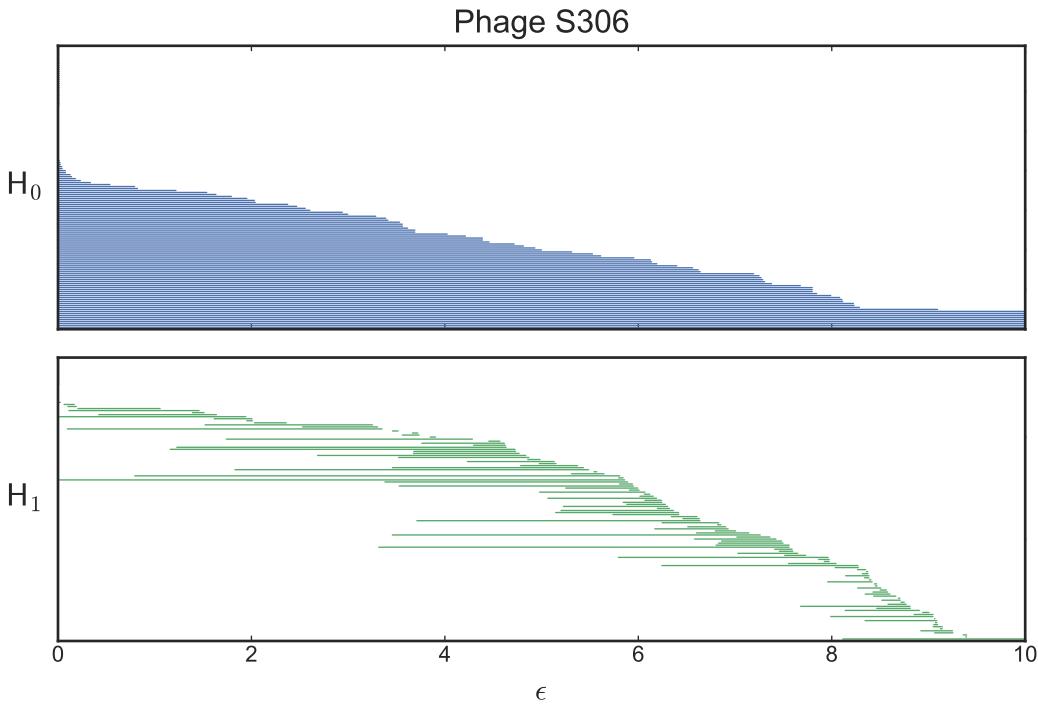


Figure 5.3: Bacteriophage Barcode Diagram using the S306 dataset

## 5.4 Representing Phage Relationships with Mapper

We used Ayasdi Mapper to construct a network representation of the phage phyletic profiles. The network was constructed using a Hamming metric on the phyletic matrix and a 2D filter function. The first filter was Metric PCA coordinate 1 with a resolution of 20 and a gain of 3.<sup>3</sup> The second filter was Metric PCA coordinate 2 with a resolution of 20 and a gain of 3. The equalize setting was used for both filter functions, which ensures that in the filtered space each bin has approximately the same number of points. This resulted in a network consisting of 201 nodes from the original 306 rows. The basic structure of the network is shown in Figure 5.6, where node color corresponds to the number of phages contained in the node. The network consists of one large connected component, two smaller connected components, and 21 singly connected nodes. The large connected component has local

---

<sup>3</sup>The parameter settings are in arbitrary units and tuned by hand to produce the most visually useful graph.

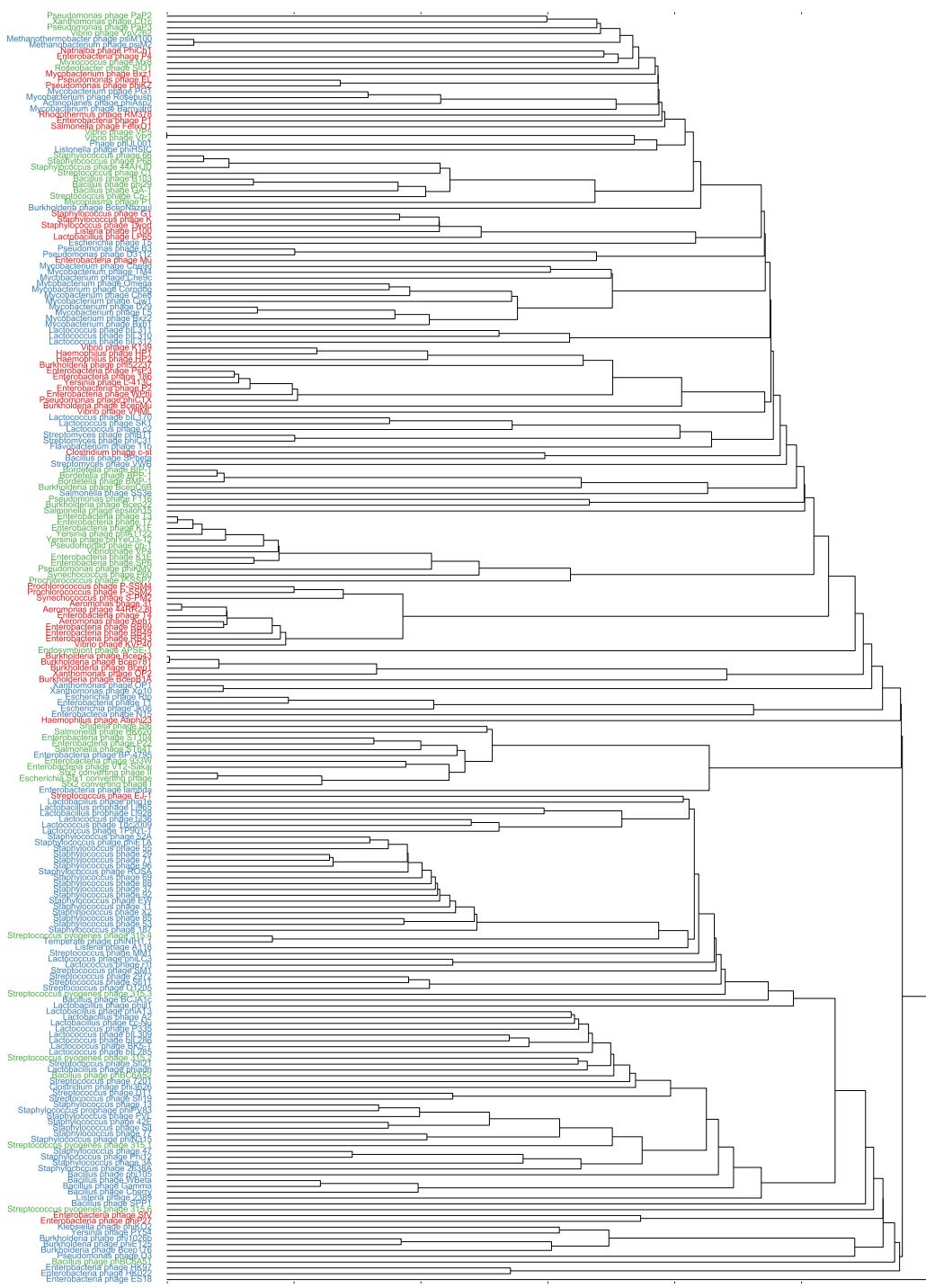


Figure 5.4: The dendrogram constructed from the  $H_0$ , restricted to bacteriophages of order Caudovirales. Myoviridae in red, Siphoviridae in blue, and Podoviridae in green. The family classifications are inconsistent with the hierarchical clustering.

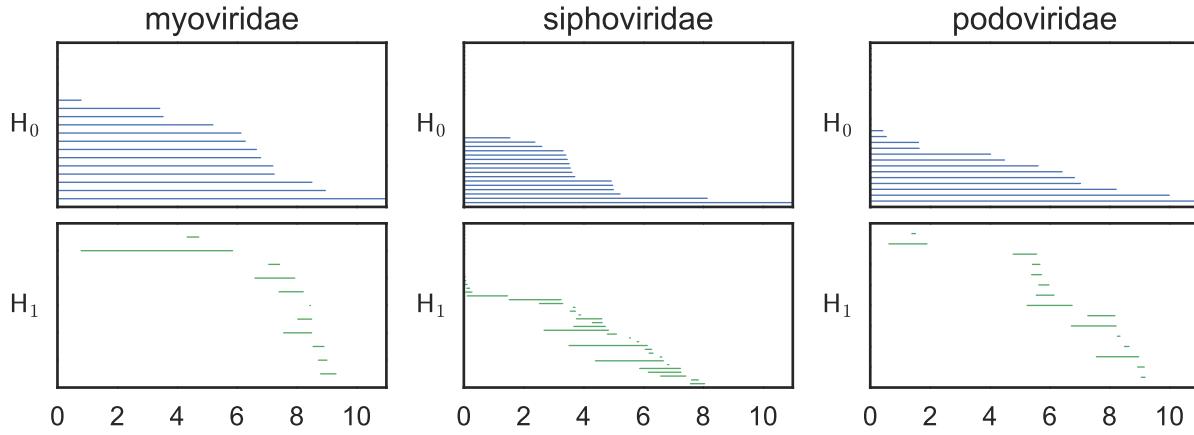


Figure 5.5: Barcode Diagrams for Families of Order Caudovirales, including Siphoviridae, Myoviridae, and Podoviridae.

regions of clustering, which will be considered further later.

We first examined how well the existing taxonomic classifications localized in the Mapper representation. If the taxonomy accurately reflects the molecular characteristics that were used to construct the network, we would expect to see strains belonging to the same level of hierarchy as localized together in the network, with minimal mixing between strains of different classification. We show the representation of the three families of order Caudovirales in Figure 5.7. Each node is colored by the proportion of rows from that family contained in the node.<sup>4</sup> We immediately see that each family is widely dispersed across the network. On closer examination, we see that the patterns of spread resemble those of the dendrogram in Figure 5.4, in that there are multiple clusters core clusters for each family. For example, the Myoviridae family has clusters in the bottom left and bottom right of the large component, and two singleton clusters. This roughly corresponds to the four clusters of Myoviridae in the  $H_0$  dendrogram.

How strongly a particular classification is reflected by a network can be quantitatively measured using a modularity score [121]. Modularity was originally devised for identify-

---

<sup>4</sup>Recall that the nodes in a Mapper network can be composed of multiple nodes, depending on the parameters of the filter function used.

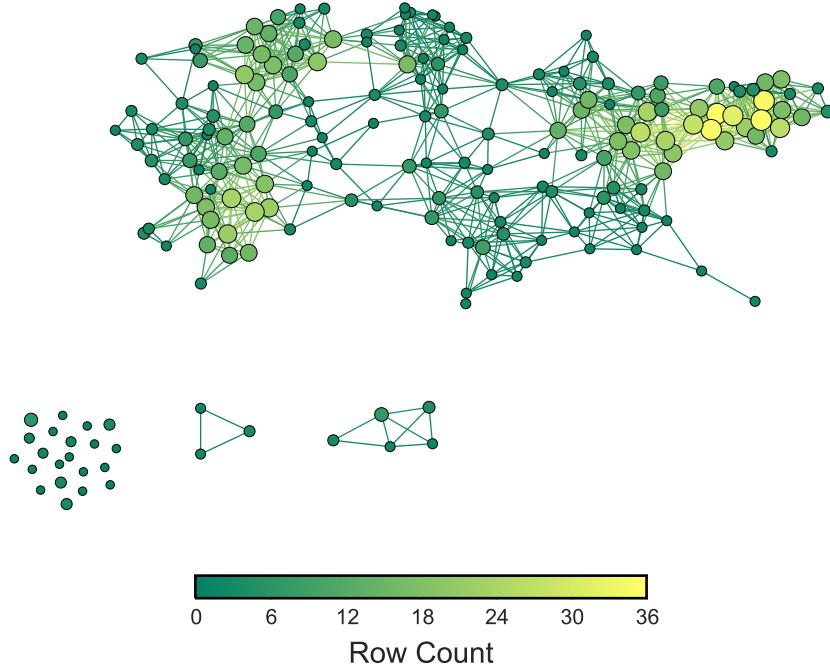


Figure 5.6: Phage Mapper Network. Network constructed using Mapper as implemented in Ayasdi Iris [6]. The network was constructed using a Hamming metric with a 2D Metric PCA filter function (resolution=20, gain=3, equalize). Nodes in the network represent clusters of phages and edges connect nodes that contain samples in common. Nodes are colored by the number of phages in each node.

ing community structure in networks. Intuitively, more tightly localized network divisions will have a higher modularity, while dispersed divisions will have a lower modularity. The standard definition for a two-class division is We use a modified form of modularity

$$Q = \frac{1}{m} \sum_{ij} A_{ij} s_i s_j \quad (5.3)$$

where  $m$  is the total number of edges in the network,  $A$  is the adjacency matrix of the network, and  $s_i = \pm 1$  is the class membership of node  $i$ .<sup>5</sup> The modularity ranges between 0 and 1. We use a strict class membership, in which  $s_i = 1$  for node  $i$  if any row in the now contains the annotation of interest. The modularity score for each family of Caudovirales

---

<sup>5</sup>The standard definition of modularity includes a term measuring how tightly connected each module is. We are only interested in the localization of each modular and neglect this term.

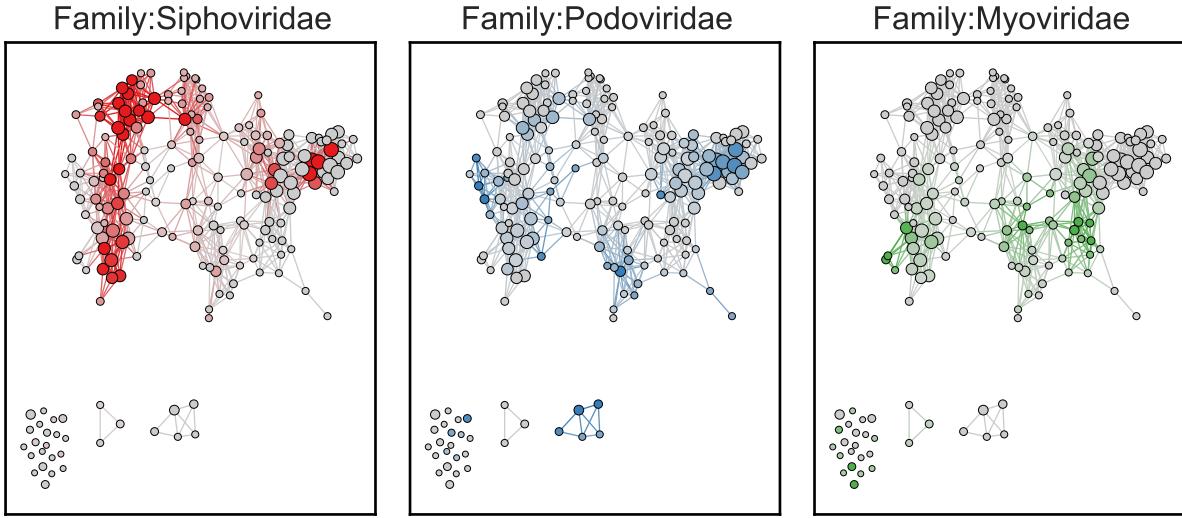


Figure 5.7: Taxonomic localization in the bacteriophage network. Network nodes are colored by presence of phages for each family in order Caudovirales.

is shown in Figure 5.8.

Second, we examined how well host correlated with network structure. We show this for the top six hosts represented in our dataset in Figure 5.9. While Enterobacteria has several pockets of representation within the network, phages are on average more strongly clustered by host than by taxonomy. This is consist with existing evidence that phages of similar host range have a common environment for reticulate exchange[101]. Modularity scores for the most dominant four hosts are shown in Figure 5.8. Staphylococcus has the highest defined modularity, which is consistent with the earlier reports about strong coupling and high levels of exchange between the Staphylococcus host and its viruses [44].

Finally, we clustered the network using the MCL graph clustering algorithm [60], as implemented in the Python MCLMarkovCluster package [99]. The MCL algorithm takes two input parameters which control the coarseness of the clustering: an expansion factor  $e$  and an inflation factor  $i$ . We set  $e = 5$  and  $i = 5$ . Ignoring the singleton nodes, this resulted in eleven clusters, as shown in Figure 5.10. For each cluster, we used a hypergeometric test to identify particular protein families that were over- or under-represented in each cluster.

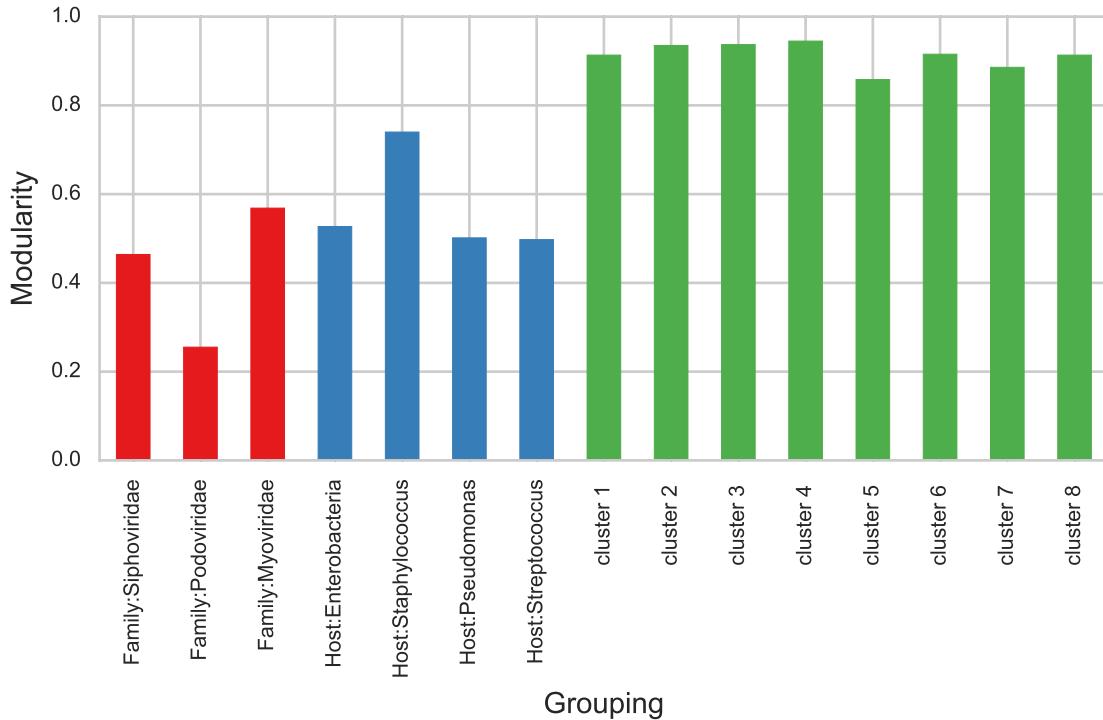


Figure 5.8: Modularity Scores for Different Divisions of the Phage Network. We show the modularities for divisions defined by taxonomic family and host range, as well as the clusters we identify using MCL.

After correcting for multiple testing, the protein families were most significantly associated with particular clusters are shown in Table 5.2.

## 5.5 Conclusions

In this chapter, we analyzed reticulate evolution in bacteriophages, using data from fully sequenced phage genomes represented as phyletic profiles measuring gene content. First, we used persistent homology to show that there are high levels of reticulate exchange across multiple taxonomic scales. Information in the  $H_0$  barcode confirmed the inconsistency of the ICTV classification. Information in the  $H_1$  barcode was used to compare levels of reticulate exchange among different phages. Second, we used Mapper to construct a network representation of phage molecular relationships. We examined how well different annotations,

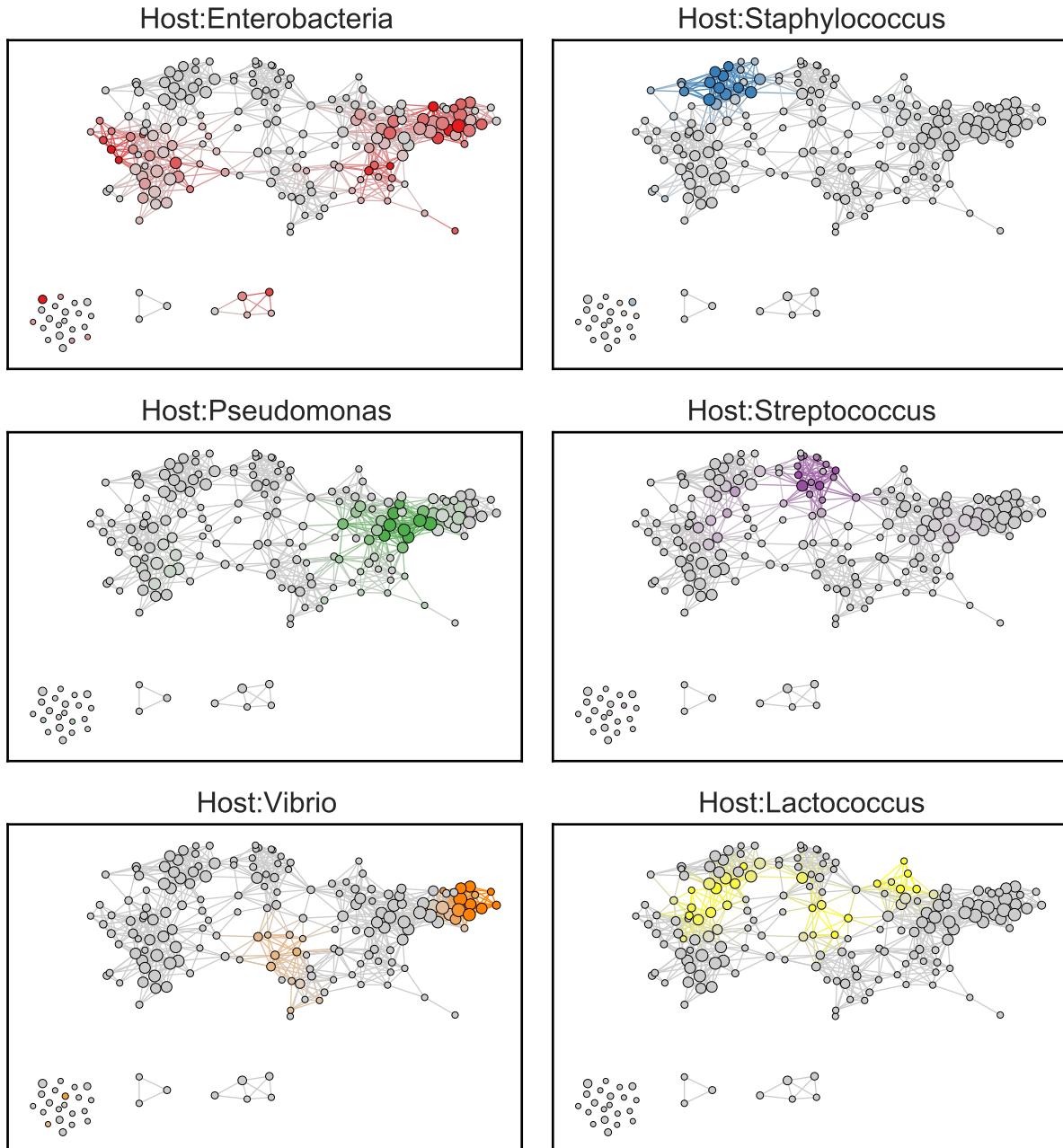


Figure 5.9: Host localization in the bacteriophage network. Compared to taxonomic family, phages are more tightly localized, reflecting the degree to which shared host range provides an environment for reticulate exchange.

Table 5.2: Phage Network MCL clustering annotations and representative protein families

Cluster	Protein Family	p-value	Function	Cluster	Protein Family	p-value	Function
Cluster 1	pf_0001	6.17e-25	tail tape measure protein	Cluster 5	pf_0008	7.45e-07	unknown
	pf_0002	6.81e-21	transcriptional repressor	Cluster 6	pf_0006	2.17e-18	NA
	pf_0003	2.69e-16	tyrosine based integrase		pf_0057	8.74e-10	unknown
	pf_0004	9.08e-12	DNA binding protein		pf_0142	3.84e-09	unknown
	pf_0008	7.37e-10	unknown		pf_0082	1.79e-08	lysis protein
	pf_0010	2.36e-08	terminase large subunit		pf_0165	2.09e-08	prohead
	pf_0006	2.87e-08	NA		pf_0188	1.12e-07	minor tail protein
	pf_0164	5.57e-08	NA		pf_0190	1.12e-07	tail
	pf_0012	1.74e-07	DNA replication initiation protein		pf_0189	1.12e-07	minor tail protein
	pf_0015	2.42e-07	scaffolding protein		pf_0204	5.81e-07	unknown
	pf_0187	2.70e-07	NA				
	pf_0217	3.35e-07	NA				
	pf_0013	4.63e-07	portal protein				
	pf_0017	8.84e-07	endolysin				
Cluster 2	pf_0131	9.38e-13	NA	Cluster 7	pf_0002	8.92e-11	transcriptional repressor
	pf_0279	2.33e-09	NA		pf_0121	2.46e-09	transcription factor
	pf_0109	1.02e-08	NA		pf_0016	3.36e-08	unknown
	pf_0434	4.38e-08	NA		pf_0122	6.52e-08	unknown
	pf_0435	4.38e-08	NA		pf_0156	3.35e-07	NA
	pf_0436	4.38e-08	NA		pf_0517	4.38e-07	NA
	pf_0010	1.84e-07	terminase large subunit		pf_0004	5.26e-07	DNA binding protein
	pf_0029	2.02e-07	major head protein		pf_0135	5.97e-07	unknown
	pf_0009	2.54e-07	NA		pf_0176	9.43e-07	post-translational regulator
	pf_0093	4.61e-07	NA		pf_0178	9.43e-07	unknown
	pf_0512	5.47e-07	NA		pf_0177	9.43e-07	transcription anti-termination protein
	pf_0017	7.12e-07	portal protein		pf_0012	9.78e-07	DNA replication initiation protein
Cluster 3					pf_0324	1.06e-06	NA
Cluster 4	pf_0049	3.44e-24	unknown	Cluster 8	pf_0497	1.24e-07	NA
	pf_0043	3.44e-24	unknown		pf_0417	8.23e-07	NA
	pf_0053	3.86e-23	unknown		pf_0002	9.70e-07	transcriptional repressor
	pf_0052	3.86e-23	unknown				
	pf_0007	6.19e-22	endolysin	Cluster 9	pf_0425	2.15e-14	NA
	pf_0058	4.39e-21	unknown		pf_0424	2.15e-14	NA
	pf_0060	4.39e-21	unknown		pf_0423	2.15e-14	NA
	pf_0061	4.39e-21	unknown		pf_0422	2.15e-14	NA
	pf_0002	2.23e-20	transcriptional repressor		pf_0421	2.15e-14	NA
	pf_0067	4.46e-20	unknown		pf_0420	2.15e-14	NA
	pf_0063	4.46e-20	unknown		pf_0146	2.15e-14	NA
	pf_0012	1.12e-19	DNA replication initiation protein		pf_0366	1.72e-13	NA
	pf_0018	2.04e-19	tail protein		pf_0316	7.74e-13	NA
	pf_0072	4.38e-19	unknown		pf_0315	7.74e-13	NA
	pf_0004	1.82e-18	DNA binding protein		pf_0273	2.58e-12	NA
	pf_0020	5.91e-18	unknown		pf_0274	2.58e-12	internal virion protein
	pf_0012	3.86e-17	unknown		pf_0138	2.58e-12	head protein
	pf_0042	5.24e-17	unknown		pf_0504	6.45e-12	NA
	pf_0001	1.54e-16	tail tape measure protein		pf_0503	6.45e-12	NA
	pf_0092	3.48e-16	unknown		pf_0183	7.09e-12	NA
					pf_0184	1.70e-11	portal protein
					pf_0185	1.70e-11	tail protein
					pf_0163	1.70e-11	tail protein
					pf_0426	4.50e-11	NA

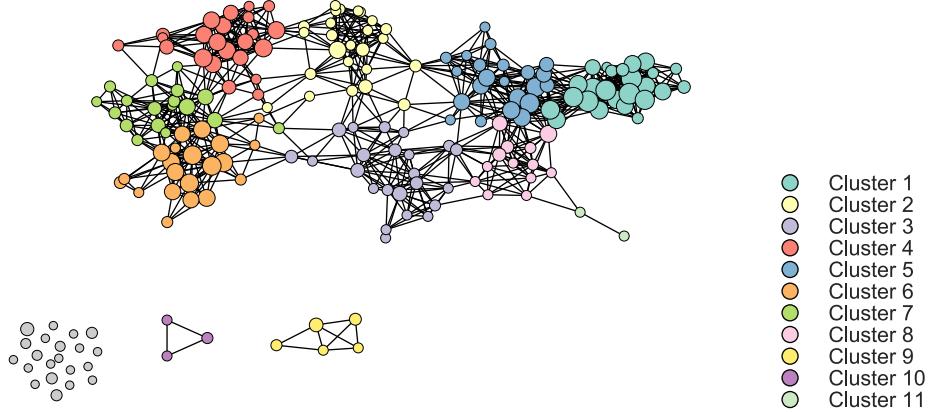


Figure 5.10: Phage Network with MCL Clustering. 11 nontrivial clusters are identified. In Table 5.2 we associate clusters with representative protein families.

including taxonomic classification and host range, localized on this network. We used a network clustering algorithm to identify communities of phages related by shared protein content, and identified protein families representative of each cluster. These clusters, while not explicitly reflecting potential phylogenetic trajectories, are more reflective of molecular similarity than existing morphological taxonomies, and can be used as a starting point for developing a more comprehensive picture of bacteriophage evolutionary dynamics. Further sequencing data will allow us to refine these clusters and provide a higher resolution



# Chapter 6

## Reassortment in Influenza Evolution

### 6.1 Introduction

In this chapter, we study influenza virus, a common human pathogen with a substantial burden on human health. Seasonal influenza epidemics have an annual mortality of between 250,000 and 500,000 [158]. Influenza pandemics, which have historically occurred roughly once every thirty years, can infect between 20-40% of the global population. For example, the Spanish influenza pandemic of 1918-1919 is estimated to have infected approximately 500 million people and lead to the death of between 50-100 million people [141]. This amounts to an infection of approximately 33% of the population and a case fatality ratio of 5-6% of global population.

The natural host reservoir of influenza is waterfowl. Within this reservoir, several distinct subtypes circulate. Subtypes are labeled by the antigenic type of two surface proteins, hemagglutinin (HA) and neuraminidase (NA).<sup>1</sup> There are presently eighteen types of HA (H1 to H18) and eleven types of NA (N1 to N11). Zoonotic adaptations have led to multiple introductions to human populations, which have resulted in both isolated outbreaks and

---

<sup>1</sup>An antigen is any molecule that elicits a host immune response. The adaptive immune system learns to recognize and protect against particular antigens. In order to evade the host immune response, the virus will mutate, giving rise to antigenic variation.

sustained transmission [118].<sup>2</sup>

The evolution of influenza is punctuated by frequent reassortment. Reassortment occurs when two virus particles coinfect the same host cell, and is a consequence of influenza having a segmented genome. The result is viral progeny that carries genomic information from two independent parental strains. This mode of evolution is known as *antigenic shift*, because it can rapidly lead to antigenically distinct viral strains.<sup>3</sup> Antigenic shifts have historically led to major pandemics, which can occur when novel surface proteins reassort with internal segments already adapted to the human host. Reassortments of this type led to Asian H2N2 flu pandemic of 1957 and the Hong Kong H3N2 flu pandemic of 1968 [105]. The 2009 H1N1 pandemic strain emerged from a triple reassortment between avian, swine, and human circulating strains [77, 135]. The pandemic had a global infection rate of between 11%-21% but a lower mortality rate than initially expected.<sup>4</sup> The 2013 H7N9 flu outbreak was caused by a triple reassortment of three distinct avian strains [35]. Traditionally, reassortments have been identified by hand, by comparing phylogenetic trees constructed from different genomic segments [119].

Recent years have seen increased concerns about the pandemic potential for zoonotic adaptation of highly pathogenic strains of influenza. Of particular concern is H5N1, which has an estimated case fatality rate of 50% (449 deaths from 846 confirmed human cases) [159], but has so far not exhibited sustained person-to-person transmission [158]. Studies in ferret models demonstrated sustained transmission in a reassortent H5N1 with as few as four mutations in the HA protein [85]. These concerns underscore the need to efficiently characterize and represent reticulate evolution in influenza. Since the 2009 H1N1 pandemic, substantial effort has been put into collecting and organizing fully sequenced influenza genomes. The

---

<sup>2</sup>Understanding the genetic basis for host adaptation is an important and controversial research area. Our work in this area in collaboration with Yoshihiro Kawaoka is forthcoming [149].

<sup>3</sup>As opposed to *antigenic drift*, due to random mutation and genetic drift.

<sup>4</sup>The 2009 H1N1 pandemic is an excellent example of the delicate balance between virulence and transmissibility.

NCBI Influenza Virus Resource now contains over 400,000 unique viral isolates [9]. The large quantity of genomic data that has been collected provides an ideal environment for studying reticulate evolution with high resolution.

## 6.2 Influenza Virology

Influenza is an enveloped single-stranded negative-sense RNA virus of family Orthomyxoviridae. The virus has a segmented genome with eight segments coding for eleven proteins. The genome length is approximately 13.5 kb. The viral structure is shown in Figure 6.1. The segments are typically ordered from longest to shortest and are detailed in Table 6.1. Of these segments, hemagglutinin (HA) and neuraminidase (NA) are the two most important. HA and NA form the two surface protein markers and are responsible for viral entry and release. HA regulates host cell binding and entry into host epithelial cells. HA is the strongest determinant of host specificity: different hosts express different sialic acid types. Avian influenza binds to type 2-3 sialic acid receptors, while human influenza binds to type 2-6 sialic acid receptors. NA is the surface protein that cleaves the newly replicated virus particles from the cell surface. Together, HA and NA determine the strain subtype and are a primary marker of host specificity and transmissibility. PA, PB1, and PB2 form a polymerase complex and are involved in viral replication. Mutations in these proteins can be among the most important in determining host adaptation and virulence, particularly mutation PB2-E627K, [139, 76]. The remaining proteins, including NP, M1, M2, and NS1 are largely structural proteins involved in capsid formation and viral packaging.

## 6.3 Influenza Reassortment

We characterized reassortment in avian influenza using persistent homology. We first compiled an aligned dataset of 3,105 complete avian influenza genomes from the NIH Influenza Sequence Database. These sequences span in time from 1956 to 2012. We collected samples

Table 6.1: Influenza Protein Segments

Segment Number	Segment Name	Protein	Length (aa)
1	Polymerase basic 2	PB2	759
2	Polymerase basic 1	PB1	757
3	Polymerase acidic	PA	716
4	Hemagglutinin	HA	563
5	Nucleoprotein	NP	498
6	Neuraminidase	NA	470
7	Matrix	M1	252
		M2	97
8	Nonstructural	NS1	230
		NS2	121

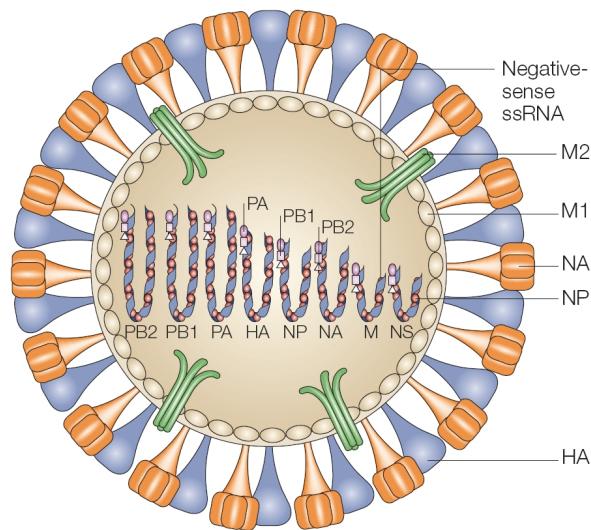


Figure 6.1: Structure of an influenza virus particle. Surface antigens HA and NA coat this surface and are involved in viral entry and exit into the host cell. The surface capsid is formed from matrix proteins M1 and M2. PB1, PB2, and PA form a polymerase complex assisting in viral replication in the infected cell.

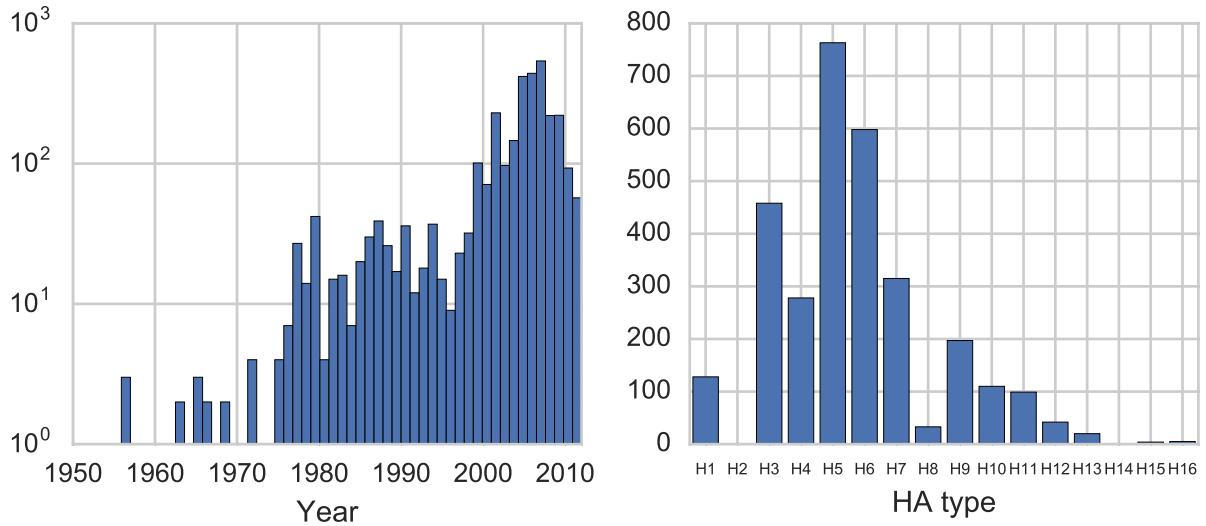


Figure 6.2: The avian influenza dataset analyzed in this chapter. Sequences spanned from 1950 to 2011, with the vast majority being collected after 2000. Most sequences were of HA type H5, with H6 and H3 following. Dataset was collected from the NCBI Influenza Virus Resource [9]

from all influenza subtypes. The majority of our sequences are of the H5 and H6 type, with a smaller proportion of H3, H7, and H9. The distribution of collected HA types and years is shown in Figure 6.2.

We first applied persistent homology to each genomic segment individually, as shown in Figure 6.3. Here we see very little higher homology, consistent with no intra-segmental recombination. The presence of higher homology is likely due to back mutation, which is expected to be more common in viruses with high mutation rates and shorter genomes (i.e. the infinite sites model does not hold). However, an analysis of the concatenated full genome reveals a complex topology, with a large number of homological invariants in one and two dimensions (Figure 6.4).

These results show that persistent homology can detect pervasive reassortment in influenza. One-dimensional ICR provides a lower-bound estimate of reassortment rate. We calculate  $\text{ICR} < 1$  event per year for classic H1N1 swine and H3N2 human influenza, supported by previous phylogenetic estimates [108, 78]. In contrast, we calculate a much higher

rate of 22.16 reassortments per year for avian influenza A. This difference could be explained by the high diversity and frequent coinfection of avian viruses [106] and correlates with the high proportion of avian reassortants reported in previous studies [54].

We used mapper to visualize the relationships in our influenza dataset. A series of mapper networks is shown in Figure 6.5. The networks were generated using a Hamming metric and the first and second MDS components as a 2D lens. In each subfigure we color the network by influenza subtype, for the top ten subtypes represented in the dataset. We can see that the current classification of flu sequences by HA and NA type is a reasonable approach, as flu isolates of the same subtype tend to cluster together tightly within the network. H6N2 is the sole subtype to be represented by two clusters in the network. When we examined the members of each H6N2 cluster, we found that both consisted of isolates spanning long time frames, suggesting that multiple stable lineages of H6N2, each carrying different internal segments, have persisted.

## 6.4 Nonrandom Association of Genome Segments

We observed nonrandom association of flu segments. Statistical inference on the loops corresponding to reassortments identified segments that tend to co-segregate with each other during reassortment. In particular, polymerases co-segregate, while genes coding for envelope and capsid proteins show independent reassortment patterns. Cosegregation of polymerases suggests that effective proteinprotein interaction between the polymerase complex and the NP protein constrain reassortment.

Although previous phylogenetic studies confirmed a high reassortment rate in avian influenza, none has identified a clear pattern of gene segment association [54]. To determine whether any segments cosegregate more than expected by chance, we considered all pairs of concatenated segments and estimated the number of reassortments using  $b_1$ . We then ascertained the significance of observing a number of reassortments between each pair of segments

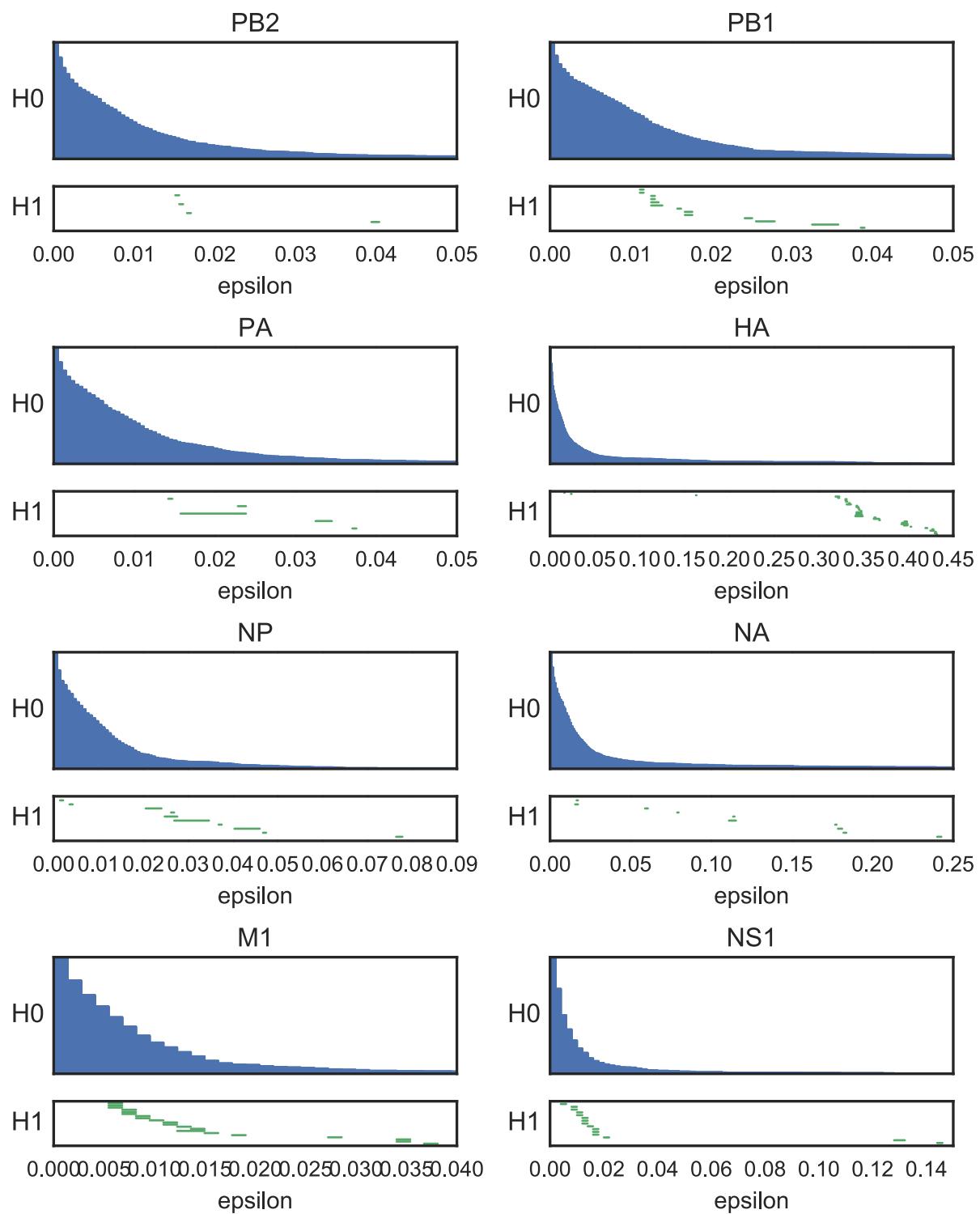


Figure 6.3: Influenza Genome Segment Barcodes. Persistent homology computed on a per-segment basis reveals very little  $H_1$  homology, indicating limited intrasegment reticulation.

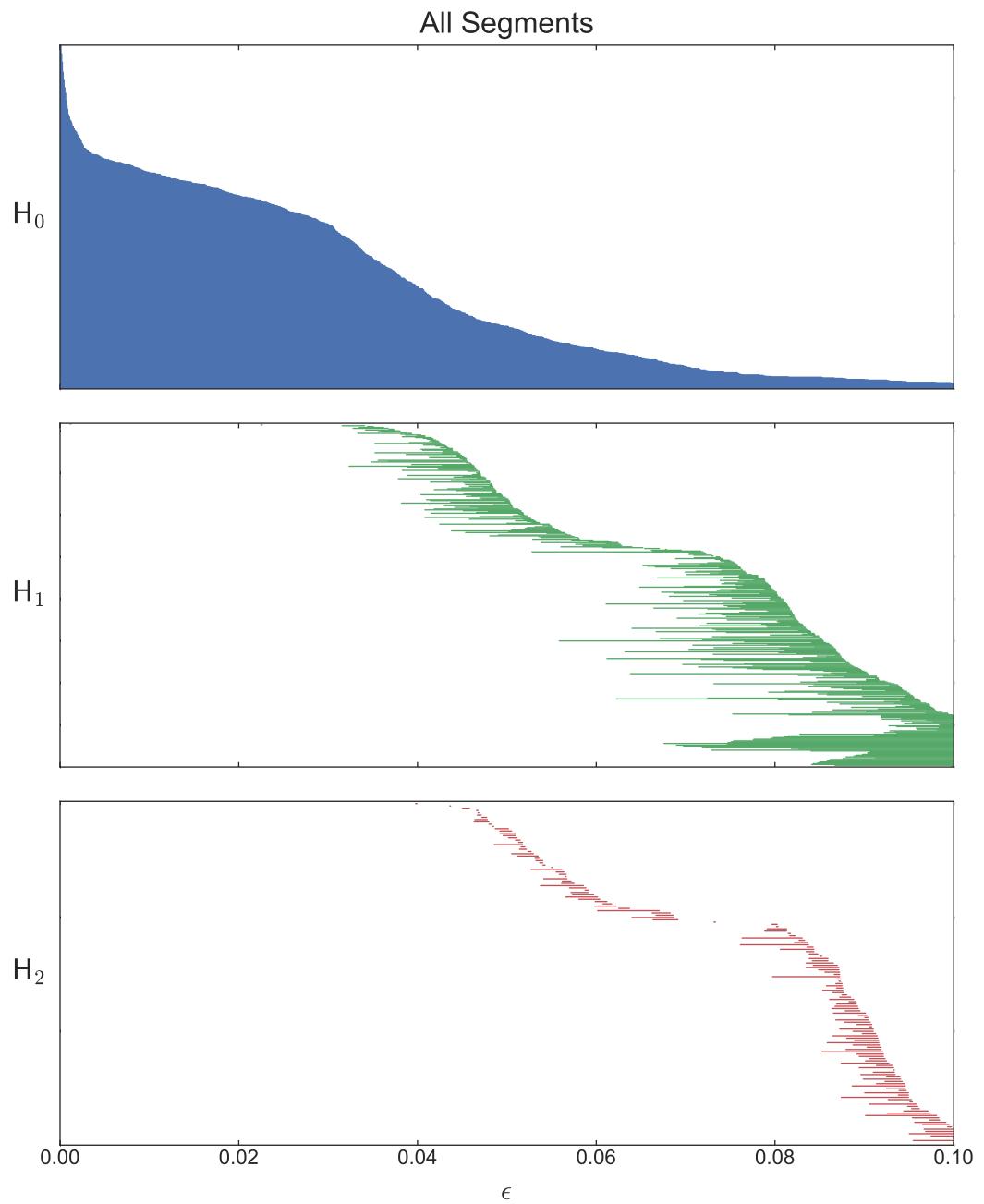


Figure 6.4: Influenza Concatenated Genome Barcode. Persistent homology computed on the full concatenated genome reveals substantial  $H_1$  and  $H_2$  homology, indicating high levels of reticulate exchange.

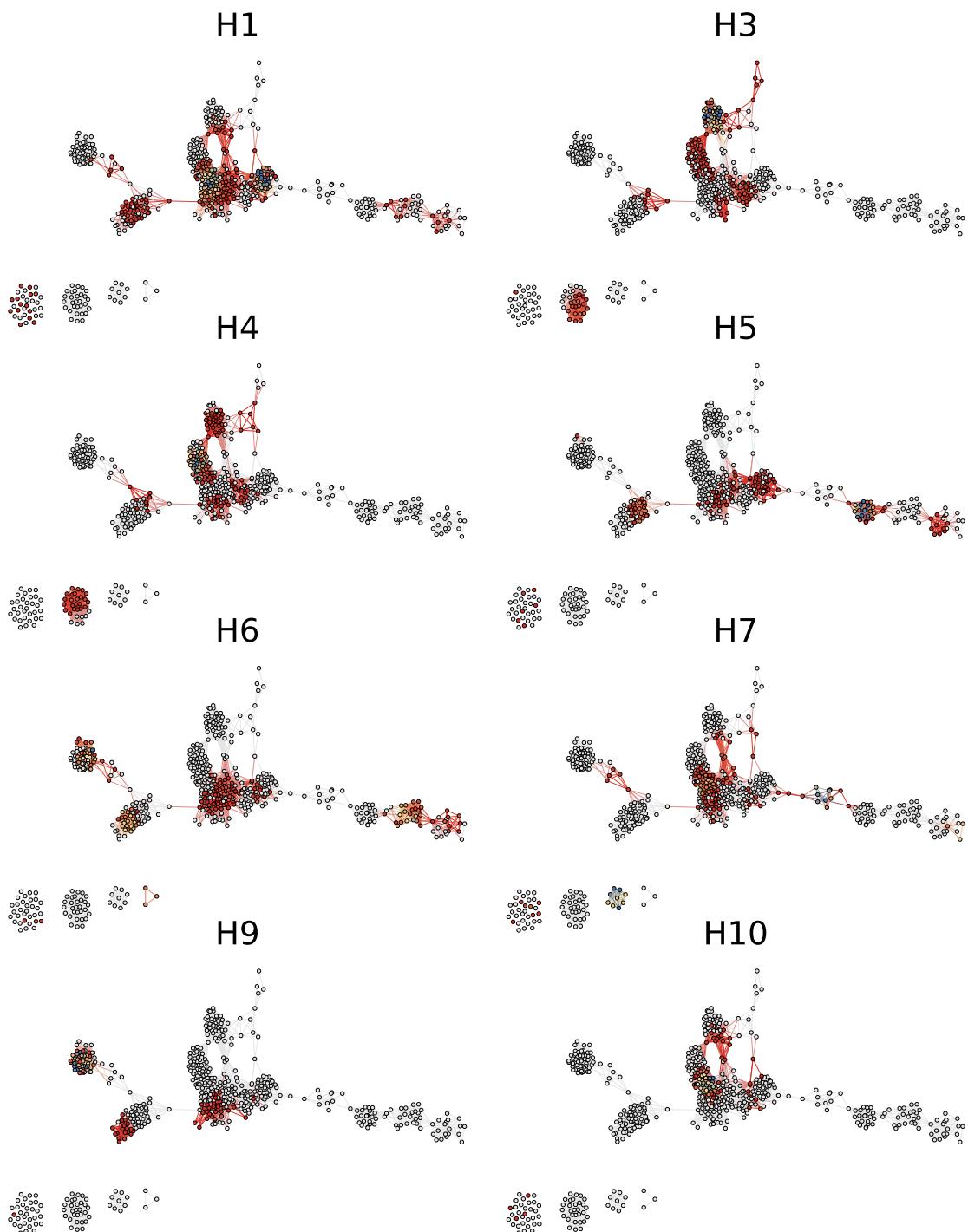


Figure 6.5: Influenza Networks By HA Subtype. The networks were generated using Ayasdi using a Hamming metric and a 2D MDS coordinate lens. Lens parameters were (gain=5, resolution=40, equalize=False).

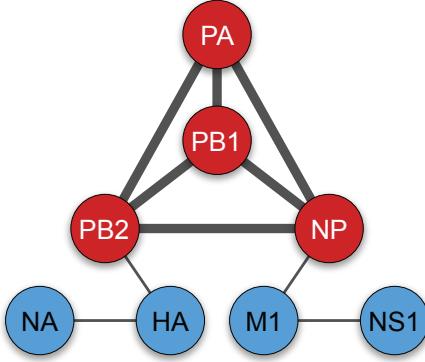


Figure 6.6: Influenza Nonrandom Reassortment

given the total estimate of reassortments in the concatenated genome. These patterns of cosegregation are represented in Figure 6.6, in which thicker edges indicate cosegregation, as measured by a lower level of homology between segment pairs. Analysis of avian influenza reveals a statistically significant configuration of four cosegregating segments: polymerase basic 2 (PB2), polymerase basic 1 (PB1), polymerase acidic (PA), and nucleoprotein (NP). Interestingly, this pattern mimics previous *in vitro* results that suggest that effective protein-protein interaction between the polymerase complex and the NP protein constrain reassortment [106].

## 6.5 Multiscale Flu Reassortment

We computed persistent homology on the avian influenza sequences across the seven major HA subtypes. The persistence diagram is shown in Figure 6.7, along with density estimates for the birth and death distributions. Both birth and death times appear strongly bimodal, unlike in the coalescent simulations, which were strictly unimodal. This suggests two distinct scales of topological structure. Using the representative cycles output by Dionysus on a subset of this data, we classified features as intrasubtype (involving one HA subtype) and intersubtype (involving multiple HA subtypes). The  $H_1$  barcode diagram for this data is

shown in the Figure 6.7 inset. Intrasubtype features, in blue, occur at an earlier filtration scale than intersubtype features, in green. The multiscale topological approach of persistent homology can distinguish biological events occurring at different genetic scales.

We isolated the two peaks and estimated two recombination rates: an intrasubtype  $\rho_1 = 9.68$ , and an intersubtype  $\rho_2 = 21.43$ . We conclude that intersubtype recombination occurs at a rate over twice that of intrasubtype recombination, however a genetic barrier exists that maintains distinct subtype populations. The nature of this barrier warrants further study. This illustrates a real-world example in which multiscale topological structure can be captured by persistent homology and given biological interpretation.

## 6.6 Conclusions

In this chapter we analyzed reassortment patterns in influenza. The segmented nature of the influenza genome, and the large amount of collected genome information, make influenza ideal for the application of topological methods. Reassortment occurs when a single cell is coinfect ed by multiple strains of the virus, and can lead to the emergence of novel pandemics. Current methods of classifying influenza, based solely on HA and NA subtype, fail to account for information in the internal segments and there is at present no consistent methodology for dealing with reassortments. We have applied methods from TDA to characterize both the scale and frequency of reassortment in influenza, estimating both reassortment rates and cosegregation patterns. Using Mapper, we determined classifications of viruses based on whole genome information that provide a higher resolution picture into extant circulating strains. Further, from the persistence diagram we identified a bimodal structure of  $H_1$  invariants, which suggests a genetic barrier maintaining subtype diversity.

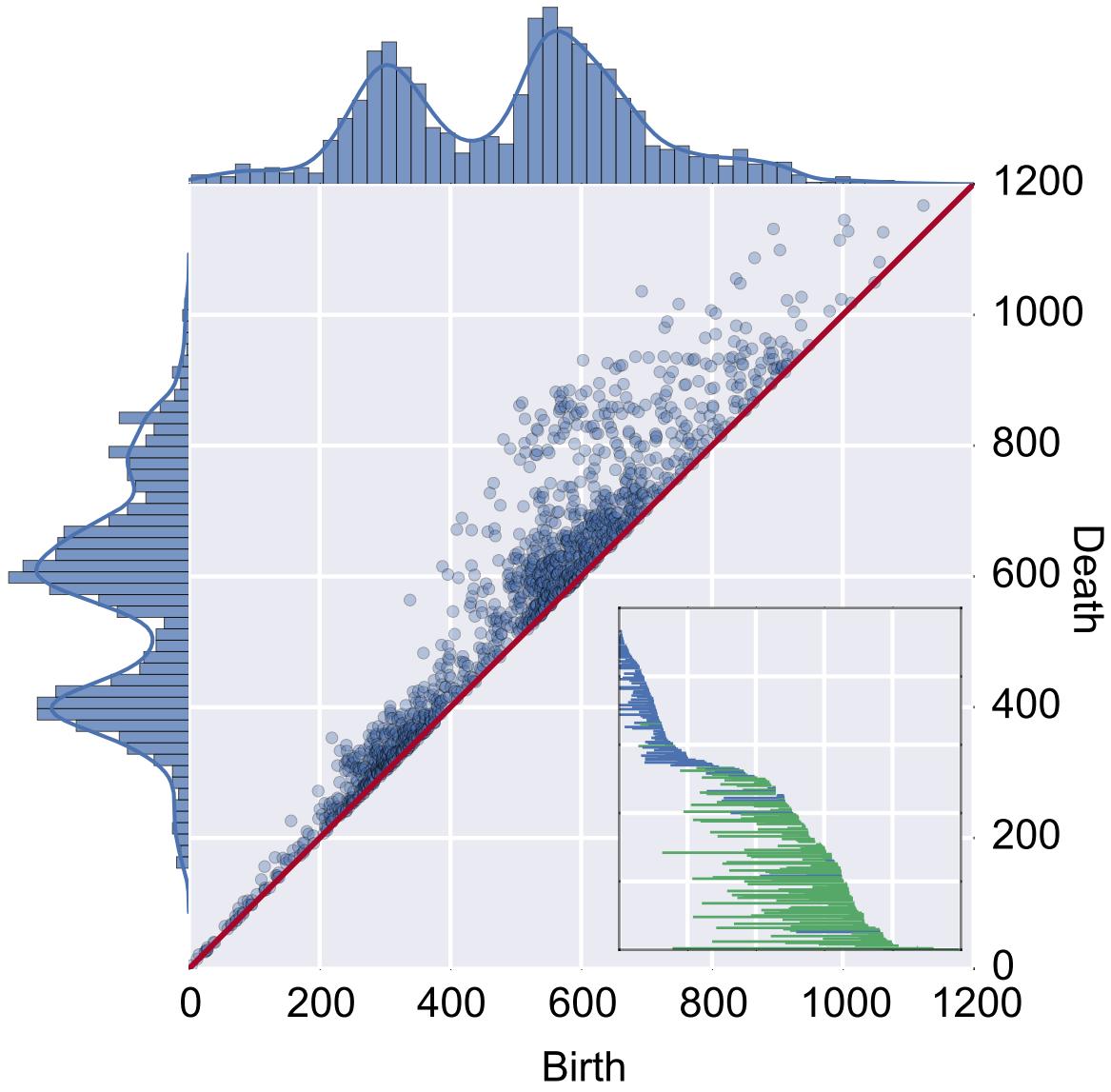


Figure 6.7: The  $H_1$  persistence diagram computed from an avian influenza dataset. On the top and left are plotted the marginal distributions of birth and death times, along with a density estimate for each distribution. The bimodality indicates two scales of topological structure. Inset: The barcode diagram for a subset of this data. Blue bars have representative cycles involving only one subtype, green bars have cycles involving multiple subtypes.

# Chapter 7

## Reticulate Evolution in Pathogenic Bacteria

### 7.1 Introduction

Pathogenic bacteria can lead to severe infection and mortality and present an enormous burden on human populations and public health systems. One of the achievements of twentieth century medicine was the development of a wide range of antibiotic drugs to control and contain the spread of pathogenic bacteria, leading to vastly increased life expectancies and global economic development. However, rapidly rising levels of multidrug antibiotic resistance in several common pathogens, including *Escherichia coli*, *Klebsiella pneumoniae*, *Staphylococcus aureus*, and *Neisseria gonorrhoeae*, is recognized as a pressing global issue with near-term consequences [120, 145, 157]. The threat of a post-antibiotic 21st century is serious, and new methods to characterize and monitor the spread of resistance are urgently needed.

Antibiotic resistance can be acquired through point mutation or through horizontal transfer of resistance genes. Horizontal exchange occurs when a donor bacteria transmits foreign DNA into a genetically distinct bacteria strain. As discussed in Chapter 2, three mechanisms

of horizontal transfer have been identified, depending on the route by which foreign DNA is acquired [123]. Foreign DNA can be acquired via uptake from an external environment (transformation), via viral-mediated processes (transduction), or via direct cell-to-cell contact between bacterial strains (conjugation). Resistance genes can be transferred between strains of the same species, or can be acquired from different species in the same environment. While the former is generally more common, an example of the latter is the phage-mediated acquisition of Shiga toxin in *E. coli* in Germany in 2011 [130]. Elements of the bacterial genome that show evidence of foreign origin are called genomic islands, and are of particular concern when associated with phenotypic effects such as virulence or antibiotic resistance.

In this chapter we explore topics relating to horizontal gene transfer in bacteria and the emergence of antibiotic resistance in pathogenic strains. We show that TDA can not only quantify gene transfer events, but also characterize the scale of gene transfer. The scale of recombination can be measured from the distribution of birth times of the  $H_1$  invariants in the barcode diagram. It has been shown that recombination rates decrease with increasing sequence divergence [67]. We characterize the rate and scale of intraspecies recombination in several pathogenic bacteria of public health concern. We select a set of pathogenic bacteria that are of public health interest based on a recently released World Health Organization (WHO) report on antimicrobial resistance [157]. Using persistent homology, we characterize the rate and scale of recombination in the core genome using multilocus sequence data. To extend our characterization to the whole genome, we use protein family annotations as a proxy for sequence composition. This allows us to compute a similarity matrix between strains. Comparing persistence diagrams gives us information about the relative scales of gene transfer at arbitrary loci. The species selected for study and the sample sizes in each analysis are specified in Table 7.1. Next, we explore the spread of antibiotic resistance genes in *S. aureus* using Mapper, an algorithm for partial clustering and visualization of high dimensional data [134]. We identify two major populations of *S. aureus*, and observe one cluster with strong enrichment for the antibiotic resistance gene *mecA*. Importantly,

Table 7.1: Pathogenic bacteria selected for study and sample sizes in each analysis.

Species	MLST profiles	PATRIC profiles
<i>Campylobacter jejuni</i>	7216	91
<i>Escherichia coli</i>	616	1621
<i>Enterococcus faecalis</i>	532	301
<i>Haemophilus influenzae</i>	1354	22
<i>Helicobacter pylori</i>	2759	366
<i>Klebsiella pneumoniae</i>	1579	161
<i>Neisseria</i> spp.	10802	234
<i>Pseudomonas aeruginosa</i>	1757	181
<i>Staphylococcus aureus</i>	2650	461
<i>Salmonella enterica</i>	1716	638
<i>Streptococcus pneumoniae</i>	9626	293
<i>Streptococcus pyogenes</i>	627	48

resistance appears to be increasingly spreading in the second population. Finally, we consider the risk of lateral transfer of resistance genes from the human microbiome into an antibiotic sensitive strain, using  $\beta$ -Lactam resistance as an example. In this environment, benign bacterial strains can harbor known resistance genes. We use a network analysis to visualize the spread of antibiotic resistance gene *mecA* into nonnative phyla. Each individual has a unique microbiome, and we speculate that microbiome typing of this sort may useful in developing personalized antibiotic therapies. These results suggest an important role for topological data mining of -omics scale data in clinical applications and personalized medicine.

## 7.2 Evolutionary Scales of Recombination in the Core Genome

Multilocus sequence typing (MLST) data was used to examine scales of recombination in the core bacterial genome. MLST is a method of rapidly assigning a sequence profile to a sample bacterial strain. For each species, a predetermined set of loci in a small number of housekeeping genes are selected as representative of the core genome of the species. At

each loci, a set of sequence types are defined by using a similarity-based clustering. As new strains are sequenced, they are annotated with a profile corresponding to the type at each locus. If a sample has a previously unseen type at a given locus, it is appended to the list of types at that locus. Large online databases have curated MLST data from labs around the world; significant pathogens can have several thousand typed strains (over 10,000 in the case of *Neisseria spp.*). Because different species will be typed at different loci, examining direct interspecies genetic exchange with this data is unfeasible, however MLST provides a large quantity of data with which to examine intraspecies exchange in the core genome. Finally, because the selected loci are primarily housekeeping genes, this type of recombination analysis will tell you only about genetic exchange in the core genome. Mobile genetic elements may have a separate rates of exchange.

We investigate genetic exchange in the twelve pathogens using MLST data from PubMLST [88]. For each strain, a pseudogenome can be constructed by concatenating the typed sequence at each locus. Using a Hamming metric, we construct a pairwise distance matrix between strains and compute persistent homology on the resulting metric space. Because of the large number of sample strains, we employ a Lazy Witness complex with 250 landmark points and  $\nu = 0$  [43]. The computation is performed using javaplex [142]. An example of our output is shown in Figure 7.1, where we plot the  $H_1$  barcode diagrams for *K. pneumoniae* and *S. enterica*. The two species have distinct recombination profiles, characterized by the range of recombinations: *K. pneumoniae* recombines at only an early short-lived scale, while *S. enterica* recombines both at the short-lived scale and a longer-lived scale. These multiple scales are reflective of population structure at the subspecies level: in *S. enterica*, there are seven defined subspecies, and experimental studies have shown high levels of reticulation both within and between subspecies groups [22]. We repeat this analysis for each species, and plot the results as a persistence diagram in Figure 7.2. Among the bulk of pathogens there appears to be three major scales of recombination, a short-lived scale at intermediate distances, a longer-lived scale at intermediate distances, and a short-lived scale at longer

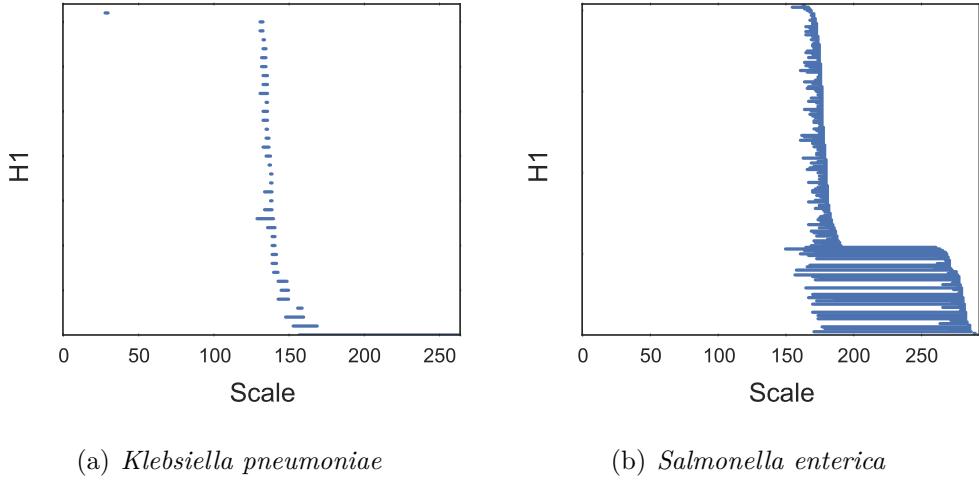


Figure 7.1: Barcode diagrams reflect different scales of core genomic exchange in *K. pneumoniae* and *S. enterica*. Differing scales can be attributed to different degrees of population substructure in the two species.

distances. We recover a pronounced pattern of reticulation at very short distances unique to *H. Pylori*, which is consistent with earlier measurements of an unusually small size of DNA imports, reported in [61].

We define a relative rate of recombination by counting the total number of  $H_1$  loops across the filtration and dividing by the number of samples for that species. The results are shown in Figure 7.3, where we observe that different species can have vastly different reticulation profiles. For example, *S. enterica* and *E. coli* have the highest reticulation rates, consistent with earlier results which have shown a high proportion of defects in the *mutS* mismatch-repair gene leading to relaxed genetic barriers to recombination [128, 102]. The low measured reticulation rate in *H. pylori* is a surprising outlier, as previous studies have shown that it lacks the mismatch repair pathways common in other bacteria, leading to higher than expected recombination rates [47]. *H. pylori* has been reported to have very little clonal structure relative to other strains, which is reflected in the star-like phylogeny that has been proposed for the species [92]. However, restriction-modification systems limiting uptake of foreign DNA have also been reported [55], suggesting that the *H. pylori* core genome is relatively resistant to reticulation at wider genomic scales. It is therefore plausible that the

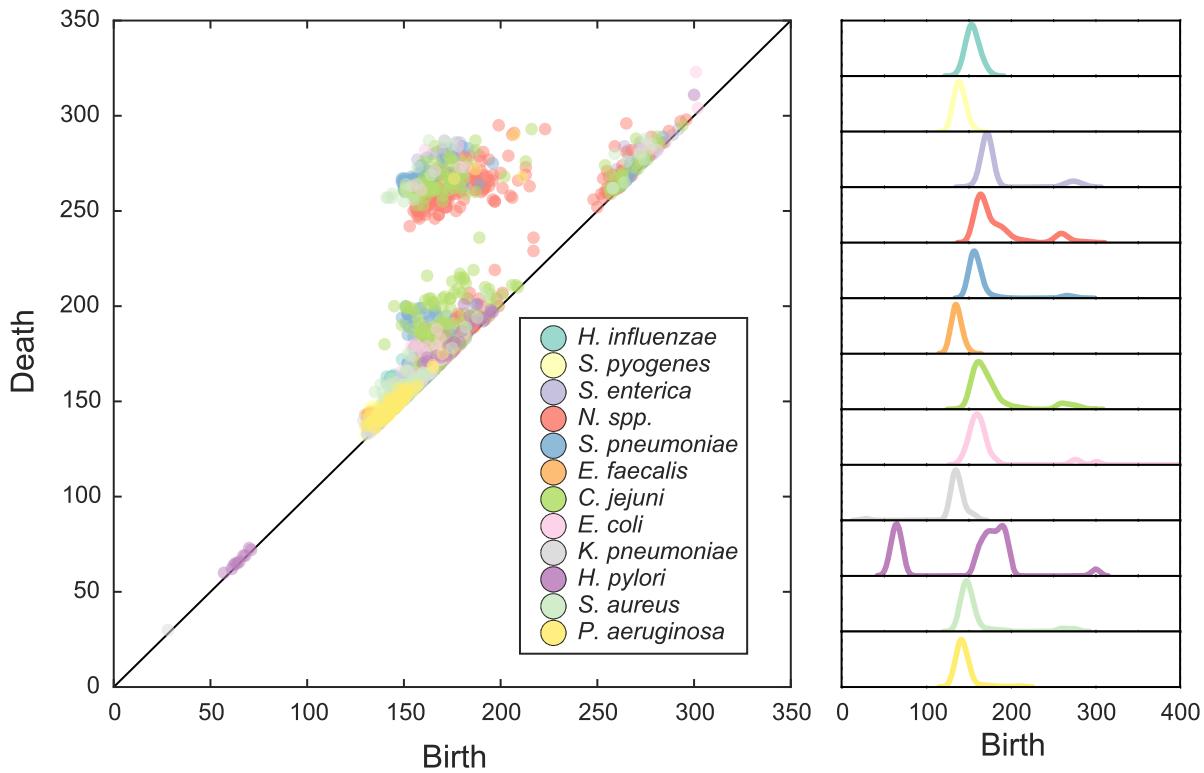


Figure 7.2: The  $H_1$  persistence diagram for the twelve pathogenic strains selected for this study using MLST profile data. There are three broad scales of recombination. To the right is the birth time distribution for each strain. *H. pylori* has an earlier scale of recombination not present in the other species, corresponding to the atypically small size of DNA imports in the species [61].

lower signal from persistent homology is due to systematic reduced sampling of particular lineages, suggesting that accounting for larger-scale population structure is important when making estimates of reticulation rates.

## 7.3 Protein Families as a Proxy for Genome Wide Reticulation

Protein family annotations cluster proteins into sets of isofunctional homologs, i.e., clusters of proteins with both similar sequence composition and similar function. A particular strain

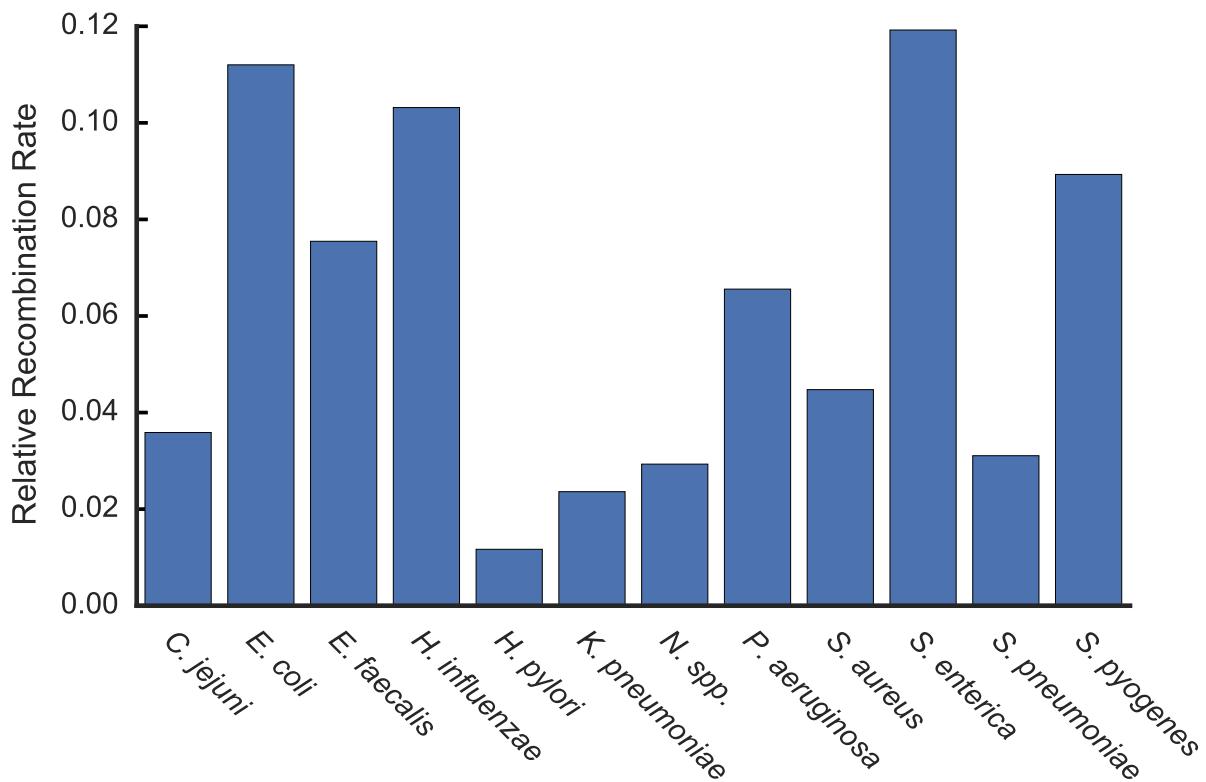


Figure 7.3: Relative recombination rates computed by persistent homology from MLST profile data.

is represented as a binary vector indicating the presence or absence of a given protein family. Correlations between strains can reveal genome-wide patterns of genetic exchange, unlike the MLST data which can only provide evidence of exchange in the core genome. We use the FigFam protein annotations in the Pathosystems Resource Institute Center (PATRIC) database because of the breadth of pathogenic strain coverage and depth of genomic annotations [151]. The FigFam annotation scheme consists of over 100,000 protein families curated from over 950,000 unique proteins [112].

For each strain we compute a transformation into FigFam space. We transform into this space because the frequency of genome rearrangements and differences in mobile genetic elements makes whole genome alignments unreliable, even for strains within the same species. As justification for performing this step, it has been shown experimentally that recombin-

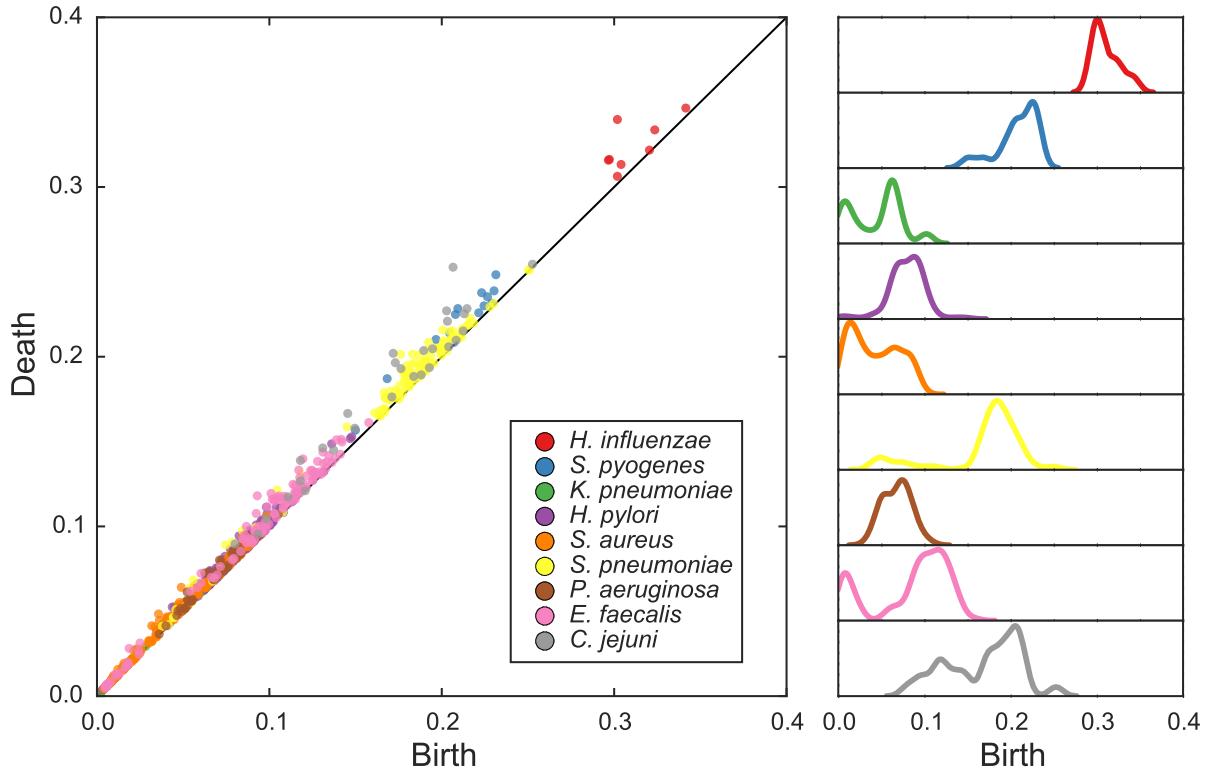


Figure 7.4: Persistence diagram for a subset of pathogenic bacteria, computed using the FigFam annotations compiled from PATRIC. Compared to the MLST persistence diagrams, the Figfam diagrams have a more diverse scale of topological structure.

tion rates decrease with increasing genetic distance [67]. After transforming, we construct a strain-strain correlation matrix and compute the persistent homology in this space. In Figure 7.4 we show the persistence diagram relating the structure and scale between different species. We find that different species have a much more diverse topological structure in this space than in MLST space, and a wide variety of recombination scales. The large scales of exchange in *H. influenzae* suggest it can regularly acquire novel genetic material from distantly related strains. This is consistent with studies that have found *H. influenzae* to be perpetually competent for uptake and integration of foreign DNA via transformation [53].

## 7.4 Antibiotic Resistance in *Staphylococcus aureus*

*S. aureus* is a gram positive bacteria commonly found in the nostrils and upper respiratory tract. Certain strains can cause severe infection in high-risk populations, particularly in the hospital setting. The emergence of antibiotic resistant *S. aureus* is therefore of significant clinical concern. Methicillin resistant *S. aureus* (MRSA) strains are resistant to  $\beta$ -lactam antibiotics including penicillin and cephalosporin. Resistance is conferred by the gene *mecA*, an element of the Staphylococcal cassette chromosome *mec* (*SCCmec*). *mecA* codes for a dysfunctional penicillin-binding protein 2a (PBP2a), which inhibits  $\beta$ -lactam antibiotic binding, the primary mechanism of action [87]. Of substantial clinical importance are methods for characterizing the spread of MRSA within the *S. aureus* population.

To address this question, we use the FigFam annotations in PATRIC, as described in the previous section. PATRIC contains genomic annotations for 461 strains of *S. aureus*, collectively spanning 3,578 protein families. We perform a clustering analysis using the Mapper algorithm as implemented in Ayasdi Iris [6]. Principal and second metric singular value decomposition are used as filter functions, with a 4x gain and an equalized resolution of 30. This results in a graph structure with two large clusters, with a smaller bridge connecting the two, as shown in Figure 7.5. The two clusters are consistent with previous phylogenetic studies using multilocus sequence data to identify two major population groups [38].

Of the 461 *S. aureus* strains in PATRIC, 142 carry the *mecA* gene. When we color nodes in the network based on an enrichment for the presence of *mecA*, we observe a much stronger enrichment in one of the two clusters. This suggests that  $\beta$ -lactam resistance has already begun to dominate in that clade, likely due to selective pressures. More strikingly, we observe that while *mecA* enrichment is not as strong in the second cluster, there is a distinct path of enrichment emanating along the connecting bridge between the two clusters and into the less enriched cluster. This suggests the hypothesis that antibiotic resistance has spread from the first cluster into the second cluster via strains intermediate to the two, and will likely continue to be selected for in the second cluster.

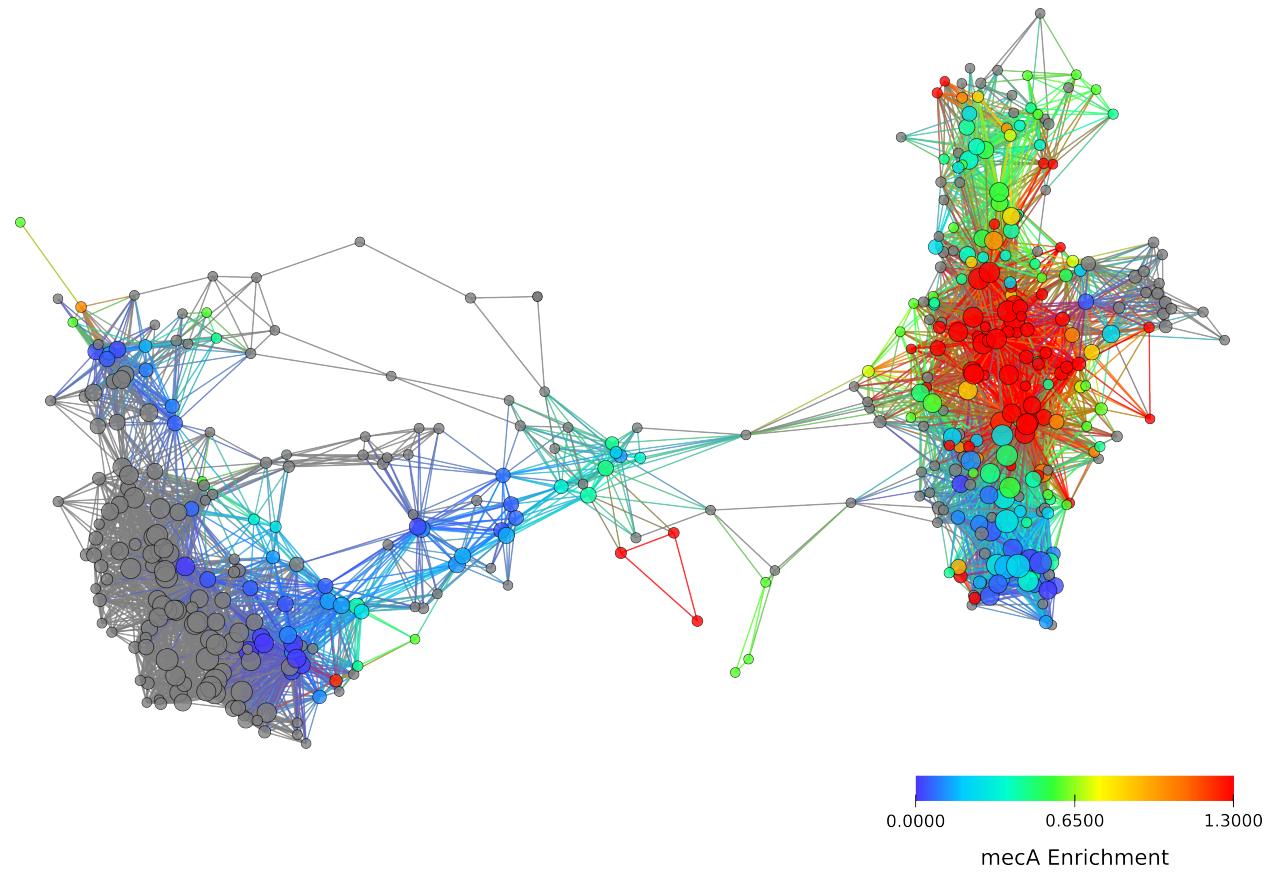


Figure 7.5: The FigFam similarity network of *S. aureus* constructed using Mapper as implemented in Ayasdi Iris [6]. We use a Hamming metric and Primary and Secondary Metric SVD filters (res: 30, gain 4x, equalize). Node color is based on strain enrichment for *meca*, the gene conferring  $\beta$ -Lactam resistance. Two distinct clades of *S. aureus* are visible, one of which has already been compromised for resistance. Of important clinical significance is the growing enrichment for *meca* in the second clade.

## 7.5 Microbiome as a Reservoir of Antibiotic Resistance Genes

While antibiotic resistance can be acquired through gene exchange between strains of the same species, it is also possible for gene exchange to occur between distantly related species. It has been recognized that an individual's microbiome, the set of microorganisms that exist symbiotically within a human host, can act as a reservoir of antimicrobial resistance genes [136, 126]. It is of substantial clinical interest to characterize to what extent an individual's

microbiome may pose a risk for a pathogenic bacteria acquiring a resistance gene through lateral transfer.

To address this question, we use data from the Human Microbiome Project (HMP), a major research initiative performing metagenomic characterization of hundreds of healthy human microbiome samples [144]. The HMP has defined a set of reference strains that have been observed in the human microbiome. We collect FigFam annotations from PATRIC for the reference strain list in the gastrointestinal tract. We focus on the gastrointestinal tract because it is an isolated environment and likely to undergo higher rates of exchange than other anatomic regions. Of the 717 reference strains, 321 had FigFam annotations. We computed a similarity matrix as in previous sections, using correlation as distance. The resulting network is shown in Figure 7.6, where strains are colored by phyla-level classifications. While largely recapitulating phylogeny, the network depicts interesting correlations between phyla, such as the loop between Firmicutes, Bacteroides, and Proteobacteria.

Next, we searched for genomic annotations relating to  $\beta$ -lactam resistance. 10 strains in the reference set had matching annotations, and we highlight those strains in the network with green diamonds. We observe resistance mostly concentrated in the Firmicutes, of which *S. aureus* is a member, however there is a strain of Proteobacteria that has acquired the resistance gene. Transfer of beta-lactam resistance into the Proteobacteria is clinically worrisome. Pathogenic Proteobacteria include *S. enterica*, *V. cholerae*, and *H. pylori*, and emergence of  $\beta$ -lactam resistance will severely impact antibiotic drug therapies.

The species composition of each individual's microbiome can differ substantially due to a wide variety of poorly understood factors [144]. In this case, an individuals personal microbiome network will differ from the network we show in Figure 7.6, which was constructed from the set of *all* strains that have been reported across studies of multiple individuals. The relative risk for acquiring self-induced resistance will therefore vary from person to person and by the infectious strain acquired. However, a network analysis of this type will give clues as to possible routes by which antibiotic resistance may be acquired. In the clinical set-

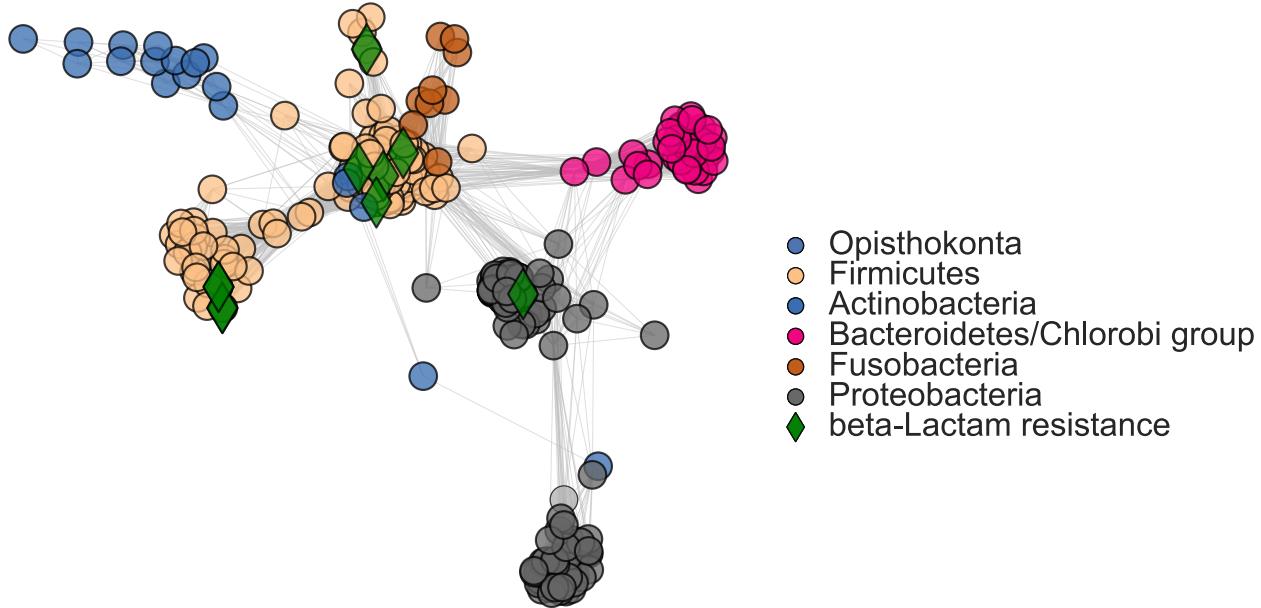


Figure 7.6: The FigFam similarity network of gastrointestinal tract reference strains identified in the Human Microbiome Project. The green diamond identifies the strains carrying resistance to  $\beta$ -Lactam antibiotics.

ting, this could assist in developing personalized antibiotic treatment regimens. We propose a more thorough expansion of this work, examining the full range of antibiotic resistance genes in order to quantify microbiome risk factors for treatment failure. We foresee an era of genomically informed infectious disease management in the clinical setting, based on an understanding of a patient’s personal microbiome profile.

## 7.6 Conclusions

In this chapter we have used some ideas from topological data analysis to bear on problems in pathogenic microbial genetics. First, we used persistent homology to evaluate recombination rates in the core genome using MLST profile data. We showed that different pathogens have different recombination rates. We expanded this to gene transfer across the whole genome by using protein family annotations in the PATRIC database. We found different scales of recombination in different pathogens. Second, we explored the spread of MRSA in *S. aureus* populations using topological methods. We noted increasing resistance in a

previously isolated population. Finally, we studied the emergence of  $\beta$ -lactam resistance in the microbiome, and proposed methods by which personal risk could be assessed by microbiome typing. These results point to a role for graph mining and topological data mining in health and personalized medicine.



# Chapter 8

## Conclusions

This thesis has considered the problem of characterizing reticulate modes of evolution in large-scale genomics data. We have drawn on methods from topological data analysis, specifically persistent homology to quantify the scale and frequency of reticulate events, and Mapper to provide condensed representations of molecular relationships. In Part I, we developed several theoretical approaches for analyzing data using TDA. In Chapter 3 we developed alternative topological complex constructions in order to increase the sensitivity of persistent homology. In Chapter 4 we developed a framework for statistical inference using the persistence diagram. We used this to develop an estimator for the recombination rate in the coalescent model, a common stochastic model in population genetics.

In Part II, we applied our general approach to several problems in evolution and genomics. In Chapter 5 we studied phages, viruses of single-celled microorganisms. We showed how persistent homology can recover inconsistencies in existing morphology-based taxonomies, used a network approach to construct an alternative genome-based representation of phage relationships, and identified representative gene families conserved within phage populations. In Chapter 6 we studied influenza, a common human pathogen. We showed how persistent homology can capture widespread patterns of reassortment, including nonrandom cosegregation of segments and barriers to subtype mixing. In contrast to traditional in-

fluenza studies, which have focused on the phylogenetic branching patterns of only the two surface-marker proteins, we used Mapper combined with whole-genome data to represent influenza molecular relationships. We identified unexpected relationships between divergent influenza subtypes. In Chapter 7 we studied pathogenic bacteria. We used two sources of data to measure rates of reticulation in both the core genome and the mobile genome across a range of species. Mapper is used to represent the population of *S. aureus* and analyze the spread of antibiotic resistance genes. The potential for the spreading of antibiotic resistance in the human microbiome is investigated.

# Bibliography

- [1] R. Adler, O. Bobrowski, M. Borman, E. Subag, and S. Weinberger, “Persistent homology for random fields and complexes,” 2010. arXiv: [1003.1001](https://arxiv.org/abs/1003.1001) [math.PR].
- [2] M. N. Alekshun and S. B. Levy, “Molecular mechanisms of antibacterial multidrug resistance,” *Cell*, vol. 128, no. 6, pp. 1037–1050, Mar. 2007. DOI: [10.1016/j.cell.2007.03.004](https://doi.org/10.1016/j.cell.2007.03.004).
- [3] S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zheng, W. Miller, and L. D. J., “Gapped BLAST and PSI-BLAST: a new generation of protein database search programs,” *Nucleic Acids Research*, vol. 25, no. 17, pp. 3389–3402, 1997. DOI: [10.1093/nar/25.17.3389](https://doi.org/10.1093/nar/25.17.3389).
- [4] M. L. Arnold, *Natural Hybridization and Evolution*, ser. Oxford Series in Ecology and Evolution. Oxford, UK: Oxford University Press, 1996.
- [5] ——, *Evolution through Genetic Exchange*. Oxford, UK: Oxford University Press, 2007.
- [6] Ayasdi Inc., *Ayasdi Core*, 2015. [Online]. Available: <http://www.ayasdi.com>.
- [7] H.-J. Bandelt and A. W. Dress, “A canonical decomposition theory for metrics on a finite set,” *Advances in Mathematics*, vol. 92, no. 1, 1992. DOI: [10.1016/0001-8708\(92\)90061-o](https://doi.org/10.1016/0001-8708(92)90061-o).
- [8] H.-J. Bandelt, P. Forster, and A. Röhl, “Median-joining networks for inferring intraspecific phylogenies,” *Molecular Biology and Evolution*, vol. 16, no. 1, pp. 37–48, 1999. DOI: [10.1093/oxfordjournals.molbev.a026036](https://doi.org/10.1093/oxfordjournals.molbev.a026036).
- [9] Y. Bao, P. Bolotov, D. Dernovoy, B. Kiryutin, L. Zaslavsky, T. Tatusova, J. Ostell, and D. Lipman, “The influenza virus resource at the National Center for Biotechnology Information,” *Journal of Virology*, vol. 82, no. 2, pp. 596–601, Jan. 2008. DOI: [10.1128/JVI.02005-07](https://doi.org/10.1128/JVI.02005-07). [Online]. Available: <http://www.ncbi.nlm.nih.gov/genomes/FLU/>.

- [10] F. Baudat, J. Buard, C. Grey, A. Fledel-Alon, C. Ober, M. Przeworski, G. Coop, and B. de Massy, “PRDM9 is a major determinant of meiotic recombination hotspots in humans and mice,” *Science*, vol. 327, no. 5967, pp. 836–840, Feb. 2010. DOI: [10.1126/science.1183439](https://doi.org/10.1126/science.1183439).
- [11] U. Bauer, M. Kerber, and J. Reininghaus, *DIPHA (a distributed persistent homology algorithm)*, version 2.1.0, 2014. [Online]. Available: <https://github.com/DIPHA/dipha/>.
- [12] ——, “Distributed computation of persistent homology,” in *Proceedings of the Sixteenth Workshop on Algorithm Engineering and Experiments (ALENEX)*, Society for Industrial & Applied Mathematics (SIAM), 2014, pp. 31–38. DOI: [10.1137/1.9781611973198.4](https://doi.org/10.1137/1.9781611973198.4).
- [13] U. Bauer, M. Kerber, J. Reininghaus, and H. Wagner, *Phat: Persistent Homology Algorithms Toolbox*, version 1.4.0, 2015. [Online]. Available: <https://bitbucket.org/phat-code/phat>.
- [14] R. Belshaw, V. Pereira, A. Katzourakis, G. Talbot, J. Paces, A. Burt, and M. Tristem, “Long-term reinfection of the human genome by endogenous retroviruses,” *Proceedings of the National Academy of Sciences*, vol. 101, no. 14, pp. 4894–4899, Apr. 2004. DOI: [10.1073/pnas.0307800101](https://doi.org/10.1073/pnas.0307800101).
- [15] L. J. Billera, S. P. Holmes, and K. Vogtmann, “Geometry of the space of phylogenetic trees,” *Advances in Applied Mathematics*, vol. 27, no. 4, pp. 733–767, 2001. DOI: [10.1006/aama.2001.0759](https://doi.org/10.1006/aama.2001.0759).
- [16] A. J. Blumberg, I. Gal, M. A. Mandell, and M. Pancia, “Robust statistics, hypothesis testing, and confidence intervals for persistent homology on metric measure spaces,” *Foundations of Computational Mathematics*, vol. 14, no. 4, pp. 745–789, 2014. DOI: [10.1007/s10208-014-9201-4](https://doi.org/10.1007/s10208-014-9201-4).
- [17] K. Borsuk, “On the imbedding of systems of compacta in simplicial complexes,” *Fundamenta Mathematicae*, vol. 35, no. 1, pp. 217–234, 1948. [Online]. Available: <http://eudml.org/doc/213158>.
- [18] L. Boto, “Horizontal gene transfer in evolution: facts and challenges,” *Proceedings of the Royal Society B: Biological Sciences*, vol. 277, no. 1683, pp. 819–827, Mar. 2010. DOI: [10.1098/rspb.2009.1679](https://doi.org/10.1098/rspb.2009.1679).
- [19] P. J. Bowler, *Evolution: The History of an Idea*. Berkeley, CA: University of California Press, 2003.

- [20] G. R. Bowman, X. Huang, Y. Yao, J. Sun, G. Carlsson, L. J. Guibas, and V. S. Pande, “Structural insight into RNA hairpin folding intermediates,” *Journal of the American Chemical Society*, vol. 130, no. 30, pp. 9676–9678, 2008. DOI: [10.1021/ja8032857](https://doi.org/10.1021/ja8032857).
- [21] L. Bren, “Bacteria-eating virus approved as food additive,” *FDA consumer*, vol. 41, no. 1, pp. 20–22, Jan. 2007.
- [22] E. W. Brown, M. K. Mammel, J. E. LeClerc, and T. A. Cebula, “Limited boundaries for extensive horizontal gene transfer among salmonella pathogens.,” *Proceedings of the National Academy of Sciences*, vol. 100, no. 26, pp. 15 676–15 681, Dec. 2003. DOI: [10.1073/pnas.2634406100](https://doi.org/10.1073/pnas.2634406100).
- [23] P. Bubenik, “Statistical topological data analysis using persistence landscapes,” *Journal of Machine Learning Research*, vol. 16, pp. 77–102, 2015. [Online]. Available: <http://www.jmlr.org/papers/v16/bubenik15a.html>.
- [24] P. Bubenik and P. T. Kim, “A statistical approach to persistent homology,” *Homology, Homotopy and Applications*, vol. 9, no. 2, pp. 337–362, 2007. DOI: [10.4310/hha.2007.v9.n2.a12](https://doi.org/10.4310/hha.2007.v9.n2.a12).
- [25] D. Burke, “Recombination in HIV: An important viral evolutionary strategy,” *Emerging Infectious Diseases*, vol. 3, no. 3, pp. 253–259, Sep. 1997. DOI: [10.3201/eid0303.970301](https://doi.org/10.3201/eid0303.970301).
- [26] P. Camara, D. Rosenbloom, K. Emmett, A. Levine, and R. Rabadan, “Fine-scale resolution of human recombination using topological data analysis,” *Cell Systems*, 2016, in press.
- [27] G. Carlsson, “Topology and data,” *Bulletin of the American Mathematical Society*, vol. 46, no. 2, pp. 255–308, 2009. DOI: [10.1090/s0273-0979-09-01249-x](https://doi.org/10.1090/s0273-0979-09-01249-x).
- [28] ——, “Topological pattern recognition for point cloud data,” *Acta Numerica*, vol. 23, pp. 289–368, 2014. DOI: [10.1017/s0962492914000051](https://doi.org/10.1017/s0962492914000051).
- [29] G. Carlsson, T. Ishkhanov, V. de Silva, and A. Zomorodian, “On the local behavior of spaces of natural images,” *International Journal of Computer Vision*, vol. 76, no. 1, pp. 1–12, 2008. DOI: [10.1007/s11263-007-0056-x](https://doi.org/10.1007/s11263-007-0056-x).
- [30] L. L. Cavalli-Sforza and A. W. Edwards, “Phylogenetic analysis. models and estimation procedures,” *American Journal of Human Genetics*, vol. 19, no. 3, pp. 550–570, 1967. DOI: [10.2307/2406616](https://doi.org/10.2307/2406616).
- [31] J. Chan, G. Carlsson, and R. Rabadan, “Topology of viral evolution,” *Proceedings of the National Academy of Sciences*, vol. 110, no. 46, pp. 18 566–18 571, Nov. 2013. DOI: [10.1073/pnas.1313480110](https://doi.org/10.1073/pnas.1313480110).

- [32] F. Chazal, D. Cohen-Steiner, L. J. Guibas, F. Mémoli, and S. Y. Oudot, “Gromov-Hausdorff stable signatures for shapes using persistence,” in *Computer Graphics Forum*, Wiley Online Library, 2009, pp. 1393–1403. DOI: [10.1111/j.1467-8659.2009.01516.x](https://doi.org/10.1111/j.1467-8659.2009.01516.x).
- [33] F. Chazal, B. T. Fasy, F. Lecci, B. Michel, A. Rinaldo, and L. Wasserman, “Sub-sampling methods for persistent homology,” in *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015, pp. 2143–2151. [Online]. Available: <http://jmlr.org/proceedings/papers/v37/chazal15.html>.
- [34] F. Chazal, M. Glisse, C. Labruére, and B. Michel, “Convergence rates for persistence diagram estimation in topological data analysis,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014, pp. 163–171. [Online]. Available: <http://jmlr.org/proceedings/papers/v32/chazal14.html>.
- [35] Y. Chen, W. Liang, S. Yang, *et al.*, “Human infections with the emerging avian influenza A H7N9 virus from wet market poultry: Clinical analysis and characterisation of viral genome,” *The Lancet*, vol. 381, no. 9881, pp. 1916–1925, Jun. 2013. DOI: [10.1016/S0140-6736\(13\)60903-4](https://doi.org/10.1016/S0140-6736(13)60903-4).
- [36] F. D. Ciccarelli, T. Doerks, C. von Mering, C. J. Creevey, B. Snel, and P. Bork, “Toward automatic reconstruction of a highly resolved tree of life,” *Science*, vol. 311, no. 5765, pp. 1283–1287, 2006. DOI: [10.1126/science.1123061](https://doi.org/10.1126/science.1123061).
- [37] G. Coop and M. Przeworski, “An evolutionary view of human recombination,” *Nature Reviews Genetics*, vol. 8, no. 1, pp. 23–34, Dec. 2006. DOI: [10.1038/nrg1947](https://doi.org/10.1038/nrg1947).
- [38] J. E. Cooper and E. J. Feil, “The phylogeny of *Staphylococcus aureus* - which genes make the best intra-species markers?” *Microbiology*, vol. 152, no. 5, pp. 1297–1305, 2006. DOI: [10.1099/mic.0.28620-0](https://doi.org/10.1099/mic.0.28620-0).
- [39] F. Crick, “Central dogma of molecular biology,” *Nature*, vol. 227, no. 5258, pp. 561–563, 1970. DOI: [10.1038/227561a0](https://doi.org/10.1038/227561a0).
- [40] T. Dagan and W. Martin, “The tree of one percent,” *Genome Biology*, vol. 7, no. 10, p. 118, 2006. DOI: [10.1186/gb-2006-7-10-118](https://doi.org/10.1186/gb-2006-7-10-118).
- [41] C. Darwin, *On the Origin of Species by Means of Natural Selection*. London, UK: Murray, 1859.
- [42] J. Davies and D. Davies, “Origins and evolution of antibiotic resistance,” *Microbiology and Molecular Biology Reviews*, vol. 74, no. 3, pp. 417–433, Aug. 2010. DOI: [10.1128/mmbr.00016-10](https://doi.org/10.1128/mmbr.00016-10).

- [43] V. de Silva and G. Carlsson, “Topological estimation using witness complexes,” in *Proceedings of the First Eurographics conference on Point-Based Graphics*, Eurographics Association, 2004, pp. 157–166. DOI: [10.2312/SPBG/SPBG04/157-166](https://doi.org/10.2312/SPBG/SPBG04/157-166).
- [44] M. Deghorain and L. Van Melderen, “The staphylococci phages family: an overview,” *Viruses*, vol. 4, no. 12, pp. 3316–3335, 2012. DOI: [10.3390/v4123316](https://doi.org/10.3390/v4123316).
- [45] W. F. Doolittle, “Phylogenetic classification and the universal tree,” *Science*, vol. 284, no. 5423, pp. 2124–2128, 1999. DOI: [10.1126/science.284.5423.2124](https://doi.org/10.1126/science.284.5423.2124).
- [46] W. F. Doolittle and R. T. Papke, “Genomics and the bacterial species problem,” *Genome Biology*, vol. 7, no. 9, p. 116, 2006. DOI: [10.1186/gb-2006-7-9-116](https://doi.org/10.1186/gb-2006-7-9-116).
- [47] M. S. Dorer, T. H. Sessler, and N. R. Salama, “Recombination and DNA repair in *Helicobacter pylori*,” *Annual Review of Microbiology*, vol. 65, no. 1, pp. 329–348, 2011. DOI: [10.1146/annurev-micro-090110-102931](https://doi.org/10.1146/annurev-micro-090110-102931).
- [48] J. W. Drake, B. Charlesworth, D. Charlesworth, and J. F. Crow, “Rates of spontaneous mutation,” *Genetics*, vol. 148, no. 4, pp. 1667–1686, Apr. 1998.
- [49] A. Dress, K. T. Huber, J. Koolen, V. Moulton, and A. Spillner, *Basic phylogenetic combinatorics*. Cambridge, UK: Cambridge University Press, 2011.
- [50] A. Dress, K. Huber, and V. Moulton, “Some variations on a theme by buneman,” *Annals of Combinatorics*, vol. 1, no. 1, pp. 339–352, 1997. DOI: [10.1007/bf02558485](https://doi.org/10.1007/bf02558485).
- [51] A. Dress, V. Moulton, and W. Terhalle, “T-theory: An overview,” *Vaccine*, vol. 17, no. 2-3, pp. 161–175, Feb. 1996. DOI: [10.1006/eujc.1996.0015](https://doi.org/10.1006/eujc.1996.0015).
- [52] A. Dress and W. Terhalle, “The tree of life and other affine buildings,” *Documenta Mathematica*, pp. 565–574, 1998. [Online]. Available: <http://eudml.org/doc/233296>.
- [53] D. Dubnau, “DNA uptake in bacteria,” *Annual Reviews in Microbiology*, vol. 53, no. 1, pp. 217–244, 1999. DOI: [10.1146/annurev.micro.53.1.217](https://doi.org/10.1146/annurev.micro.53.1.217).
- [54] V. G. Dugan, R. Chen, D. J. Spiro, *et al.*, “The evolutionary genetics and emergence of avian influenza viruses in wild birds,” *PLoS Pathogens*, vol. 4, no. 5, e1000076, May 2008. DOI: [10.1371/journal.ppat.1000076](https://doi.org/10.1371/journal.ppat.1000076).
- [55] G. R. Dwivedi, E. Sharma, and D. N. Rao, “*Helicobacter pylori* DprA alleviates restriction barrier for incoming DNA,” *Nucleic Acids Research*, vol. 41, no. 5, pp. 3274–3288, Mar. 2013. DOI: [10.1093/nar/gkt024](https://doi.org/10.1093/nar/gkt024).

- [56] H. Edelsbrunner and J. Harer, *Computational Topology: An Introduction*. American Mathematical Society, 2010. DOI: [10.1090/mhk/069](https://doi.org/10.1090/mhk/069).
- [57] K. J. Emmett and R. Rabadan, “Characterizing scales of genetic recombination and antibiotic resistance in pathogenic bacteria using topological data analysis,” in *Brain Informatics and Health*, ser. Lecture Notes in Computer Science, D. Slezak, A.-H. Tan, J. F. Peters, and L. Schwabe, Eds., vol. 8609, Springer, 2014, pp. 540–551. DOI: [10.1007/978-3-319-09891-3\\_49](https://doi.org/10.1007/978-3-319-09891-3_49).
- [58] K. Emmett and R. Rabadan, “Quantifying reticulation in phylogenetic complexes using homology,” in *BICT 2015 Special Track on Topology-driven bio-inspired methods and models for complex systems (TOPDRIM4BIO)*, 2015.
- [59] K. Emmett, D. Rosenbloom, P. Camara, and R. Rabadan, “Parametric inference using persistence diagrams: A case study in population genetics,” in *ICML Workshop on Topological Methods in Machine Learning*, 2014.
- [60] A. J. Enright, S. Van Dongen, and C. A. Ouzounis, “An efficient algorithm for large-scale detection of protein families,” *Nucleic Acids Research*, vol. 30, no. 7, pp. 1575–1584, Apr. 2002. DOI: [10.1093/nar/30.7.1575](https://doi.org/10.1093/nar/30.7.1575).
- [61] D. Falush, C. Kraft, N. S. Taylor, P. Correa, J. G. Fox, M. Achtman, and S. Suerbaum, “Recombination and mutation during long-term gastric colonization by *Helicobacter pylori*: Estimates of clock rates, recombination size, and minimal age,” *Proceedings of the National Academy of Sciences*, vol. 98, no. 26, pp. 15 056–15 061, Dec. 2001. DOI: [10.1073/pnas.251396098](https://doi.org/10.1073/pnas.251396098).
- [62] B. T. Fasy, F. Lecci, A. Rinaldo, L. Wasserman, S. Balakrishnan, and A. Singh, “Confidence sets for persistence diagrams,” *Ann. Statist.*, vol. 42, no. 6, pp. 2301–2339, DOI: [10.1214/14-AOS1252](https://doi.org/10.1214/14-AOS1252).
- [63] B. Fasy, J. Kim, F. Lecci, C. Maria, and V. Rouvreau, *Tda: Statistical tools for topological data analysis*, version 1.4.1, 2015. [Online]. Available: <https://cran.r-project.org/web/packages/TDA/index.html>.
- [64] J. Felsenstein, *Inferring Phylogenies*. Sunderland, Massachusetts: Sinauer Associates, 2004, ISBN: 978-0-87893-177-4.
- [65] W. M. Fitch and E. Margoliash, “Construction of phylogenetic trees,” *Science*, vol. 155, no. 3760, pp. 279–284, 1967. DOI: [10.1126/science.155.3760.279](https://doi.org/10.1126/science.155.3760.279).
- [66] L. R. Foulds and R. L. Graham, “The Steiner problem in phylogeny is NP-complete,” *Advances in Applied Mathematics*, vol. 3, no. 1, pp. 43–49, Mar. 1982. DOI: [10.1016/s0196-8858\(82\)80004-3](https://doi.org/10.1016/s0196-8858(82)80004-3).

- [67] C. Fraser, W. P. Hanage, and B. G. Spratt, “Recombination and the nature of bacterial speciation,” *Science*, vol. 315, no. 5811, pp. 476–480, 2007. DOI: [10.1126/science.1127573](https://doi.org/10.1126/science.1127573).
- [68] B. Gärtner, “Fast and robust smallest enclosing balls,” in *Algorithms-ESA 99*, Springer, 1999, pp. 325–338. DOI: [10.1007/3-540-48481-7\\_29](https://doi.org/10.1007/3-540-48481-7_29).
- [69] R. Ghrist, “Barcodes: The persistent topology of data,” *Bulletin of the American Mathematical Society*, vol. 45, no. 01, pp. 61–76, 2007. DOI: [10.1090/s0273-0979-07-01191-3](https://doi.org/10.1090/s0273-0979-07-01191-3).
- [70] J. P. Gogarten, W. F. Doolittle, and J. G. Lawrence, “Prokaryotic evolution in light of gene transfer,” *Molecular Biology and Evolution*, vol. 19, no. 12, pp. 2226–2238, 2002. DOI: [10.1093/oxfordjournals.molbev.a004046](https://doi.org/10.1093/oxfordjournals.molbev.a004046).
- [71] J. P. Gogarten and J. P. Townsend, “Horizontal gene transfer, genome innovation and evolution,” *Nature*, vol. 3, no. 9, pp. 679–687, 2005. DOI: [10.1038/nrmicro1204](https://doi.org/10.1038/nrmicro1204).
- [72] N. Goldenfeld and C. Woese, “Biology’s next revolution,” *Nature*, vol. 445, no. 7126, pp. 369–369, Jan. 2007. DOI: [10.1038/445369a](https://doi.org/10.1038/445369a).
- [73] S. J. Gould, *The Structure of Evolutionary Theory*. Cambridge, MA: Harvard University Press, 2002.
- [74] M. Gromov, “Hyperbolic groups,” English, in *Essays in Group Theory*, ser. Mathematical Sciences Research Institute Publications, S. Gersten, Ed., vol. 8, Springer, 1987, pp. 75–263. DOI: [10.1007/978-1-4613-9586-7\\_3](https://doi.org/10.1007/978-1-4613-9586-7_3).
- [75] A. Hatcher, *Algebraic Topology*. Cambridge, UK: Cambridge University Press, 2002.
- [76] M. Hatta, P. Gao, P. Halfmann, and Y. Kawaoka, “Molecular basis for high virulence of hong kong H5N1 influenza A viruses,” *Science*, vol. 293, no. 5536, pp. 1840–1842, Sep. 2001. DOI: [10.1126/science.1062882](https://doi.org/10.1126/science.1062882).
- [77] C. X. Hernandez, J. M. Chan, H. Khiabanian, and R. Rabidan, “Understanding the origins of a pandemic virus,” 2011. arXiv: [1104.4568v1 \[q-bio.PE\]](https://arxiv.org/abs/1104.4568v1).
- [78] E. C. Holmes, E. Ghedin, N. Miller, *et al.*, “Whole-genome analysis of human influenza A virus reveals multiple persistent lineages and reassortment among recent H3N2 viruses,” *PLoS Biology*, vol. 3, no. 9, e300, Jul. 2005. DOI: [10.1371/journal.pbio.0030300](https://doi.org/10.1371/journal.pbio.0030300).
- [79] K. T. Huber, V. Moulton, P. Lockhart, and A. Dress, “Pruned median networks: a technique for reducing the complexity of median networks,” *Molecular Phylogenetics and Evolution*, vol. 19, no. 2, pp. 302–310, 2001. DOI: [10.1006/mpev.2001.0935](https://doi.org/10.1006/mpev.2001.0935).

- [80] R. R. Hudson, “Estimating the recombination parameter of a finite population model without selection,” *Genetical research*, vol. 50, no. 3, pp. 245–250, Dec. 1987. DOI: [10.1017/s0016672300023776](https://doi.org/10.1017/s0016672300023776).
- [81] ——, “Generating samples under a Wright–Fisher neutral model of genetic variation,” *Bioinformatics*, vol. 18, no. 2, 2002. DOI: [10.1093/bioinformatics/18.2.337](https://doi.org/10.1093/bioinformatics/18.2.337).
- [82] 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*, vol. 111, no. 1, pp. 147–164, 1985. [Online]. Available: <http://genetics.org/content/111/1/147>.
- [83] D. H. Huson, R. Rupp, and C. Scornavacca, *Phylogenetic Networks: Concepts, Algorithms and Applications*. Cambridge, UK: Cambridge University Press, 2010.
- [84] J. Huxley, *Evolution: The Modern Synthesis*. Cambridge, MA: MIT Press, 1942.
- [85] M. Imai, T. Watanabe, M. Hatta, *et al.*, “Experimental adaptation of an influenza H5 HA confers respiratory droplet transmission to a reassortant H5 HA/H1N1 virus in ferrets,” *Nature*, vol. 486, no. 7403, pp. 420–428, May 2012. DOI: [doi:10.1038/nature10831](https://doi.org/10.1038/nature10831).
- [86] International Committee on Taxonomy of Viruses, *Virus Taxonomy: Ninth Report of the International Committee on Taxonomy of Viruses*, A. M. Q. King, M. J. Adams, E. B. Carstens, and E. J. Lefkowitz, Eds., ser. Immunology and Microbiology 2011. Academic Press, 2012.
- [87] S. O. Jensen and B. R. Lyon, “Genetics of antimicrobial resistance in *Staphylococcus aureus*,” *Future Microbiology*, vol. 4, no. 5, pp. 565–582, 2009. DOI: [10.2217/fmb.09.30](https://doi.org/10.2217/fmb.09.30).
- [88] K. A. Jolley and M. C. Maiden, “Bigsdb: Scalable analysis of bacterial genome variation at the population level,” *BMC Bioinformatics*, vol. 11, no. 1, p. 595, 2010. DOI: [10.1186/1471-2105-11-595](https://doi.org/10.1186/1471-2105-11-595).
- [89] T. Kaczynski, K. Mischaikow, and M. Mrozek, *Computational Homology*, ser. Applied Mathematical Sciences. Springer, Jan. 2004, vol. 157.
- [90] M. Kahle, “Random geometric complexes,” *Discrete & Computational Geometry*, vol. 45, no. 3, pp. 553–573, 2011. DOI: [10.1007/s00454-010-9319-3](https://doi.org/10.1007/s00454-010-9319-3).
- [91] ——, “Topology of random simplicial complexes: A survey,” 2013. arXiv: [1301 . 7165v1 \[math.AT\]](https://arxiv.org/abs/1301.7165v1).
- [92] A. Kalia, A. K. Mukhopadhyay, G. Dailide, Y. Ito, T. Azuma, B. C. Y. Wong, and D. E. Berg, “Evolutionary dynamics of insertion sequences in *Helicobacter pylori*,”

*Journal of bacteriology*, vol. 186, no. 22, pp. 7508–7520, Oct. 2004. DOI: [10.1128/JB.186.22.7508-7520.2004](https://doi.org/10.1128/JB.186.22.7508-7520.2004).

- [93] E. C. Keen, “Phage therapy: Concept to cure,” *Frontiers in microbiology*, vol. 3, Jul. 2012. DOI: [10.3389/fmicb.2012.00238](https://doi.org/10.3389/fmicb.2012.00238).
- [94] E. V. Koonin, “Darwinian evolution in the light of genomics,” *Nucleic Acids Research*, vol. 37, no. 4, pp. 1011–1034, Dec. 2008. DOI: [10.1093/nar/gkp089](https://doi.org/10.1093/nar/gkp089).
- [95] E. V. Koonin and Y. I. Wolf, “Genomics of bacteria and archaea: The emerging dynamic view of the prokaryotic world,” *Nucleic Acids Research*, vol. 36, no. 21, pp. 6688–6719, 2008. DOI: [10.1093/nar/gkn668](https://doi.org/10.1093/nar/gkn668).
- [96] M. Kreitman, “Nucleotide polymorphism at the alcohol dehydrogenase locus of *Drosophila melanogaster*,” *Nature*, vol. 304, no. 5925, pp. 412–417, Aug. 1983. DOI: [10.1038/304412a0](https://doi.org/10.1038/304412a0).
- [97] J. B. Kruskal, “Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis,” *Psychometrika*, vol. 29, no. 1, pp. 1–27, 1964. DOI: [10.1007/bf02289565](https://doi.org/10.1007/bf02289565).
- [98] R. Kwitt, S. Huber, M. Niethammer, W. Lin, and U. Bauer, “Statistical topological data analysis - a kernel perspective,” in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, Eds., Curran Associates, Inc., 2015, pp. 3052–3060. [Online]. Available: <http://papers.nips.cc/paper/5887-statistical-topological-data-analysis-a-kernel-perspective>.
- [99] G. Lami, *Mcl markov cluster*, version 0.3, 2014. [Online]. Available: [https://github.com/koteth/python\\_mcl](https://github.com/koteth/python_mcl).
- [100] E. S. Lander, L. M. Linton, B. Birren, *et al.*, “Initial sequencing and analysis of the human genome,” *Nature*, vol. 409, no. 6822, pp. 860–921, Feb. 2001. DOI: [10.1038/35057062](https://doi.org/10.1038/35057062).
- [101] J. G. Lawrence, G. F. Hatfull, and R. W. Hendrix, “Imbroglios of viral taxonomy: Genetic exchange and failings of phenetic approaches,” *Journal of Bacteriology*, vol. 184, no. 17, pp. 4891–4905, 2002. DOI: [10.1128/jb.184.17.4891-4905.2002](https://doi.org/10.1128/jb.184.17.4891-4905.2002).
- [102] J. E. LeClerc, B. Li, W. L. Payne, and T. A. Cebula, “High mutation frequencies among *Escherichia coli* and *Salmonella Pathogens*,” *Science*, vol. 274, no. 5290, pp. 1208–1211, Nov. 1996. DOI: [10.1126/science.274.5290.1208](https://doi.org/10.1126/science.274.5290.1208).
- [103] L. Li, W.-Y. Cheng, B. S. Glicksberg, O. Gottesman, R. Tamler, R. Chen, E. P. Bottinger, and J. T. Dudley, “Identification of type 2 diabetes subgroups through

topological analysis of patient similarity,” *Science Translational Medicine*, vol. 7, no. 311, 311ra174, Oct. 2015. DOI: [10.1126/scitranslmed.aaa9364](https://doi.org/10.1126/scitranslmed.aaa9364).

- [104] G. Lima-Mendez, J. Van Helden, A. Toussaint, and R. Leplae, “Reticulate representation of evolutionary and functional relationships between phage genomes,” *Molecular Biology and Evolution*, vol. 25, no. 4, pp. 762–777, 2008. DOI: [10.1093/molbev/msn023](https://doi.org/10.1093/molbev/msn023).
- [105] S. E. Lindstrom, N. J. Cox, and A. Klimov, “Genetic analysis of human H2N2 and early H3N2 influenza viruses, 1957–1972: Evidence for genetic divergence and multiple reassortment events,” *Virology*, vol. 328, no. 1, pp. 101–119, Oct. 2004. DOI: [10.1016/j.virol.2004.06.009](https://doi.org/10.1016/j.virol.2004.06.009).
- [106] M. D. Lubeck, P. Palese, and J. L. Schulman, “Nonrandom association of parental genes in influenza A virus recombinants,” *Virology*, vol. 95, no. 1, pp. 269–274, 1979. DOI: [10.1016/0042-6822\(79\)90430-6](https://doi.org/10.1016/0042-6822(79)90430-6).
- [107] P. Y. Lum, G. Singh, A. Lehman, T. Ishkanov, M. Vejdemo-Johansson, M. Alagappan, J. Carlsson, and G. Carlsson, “Extracting insights from the shape of complex data using topology,” *Scientific Reports*, vol. 3, Feb. 2013. DOI: [10.1038/srep01236](https://doi.org/10.1038/srep01236).
- [108] S. J. Lycett, G. Baillie, E. Coulter, *et al.*, “Estimating reassortment rates in co-circulating eurasian swine influenza viruses,” *Journal of General Virology*, vol. 93, no. 11, pp. 2326–2336, Nov. 2012. DOI: [10.1099/vir.0.044503-0](https://doi.org/10.1099/vir.0.044503-0).
- [109] R. MacPherson and B. Schweinhart, “Measuring shape with topology,” *Journal of Mathematical Physics*, vol. 53, no. 7, p. 073516, 2012. DOI: [10.1063/1.4737391](https://doi.org/10.1063/1.4737391).
- [110] W. P. Maddison, “Gene trees in species trees,” *Systematic Biology*, vol. 46, no. 3, pp. 523–536, Sep. 1997. DOI: [10.2307/2413694](https://doi.org/10.2307/2413694).
- [111] C. Maria, J.-D. Boissonnat, M. Glisse, and M. Yvinec, “The GUDHI library: Simplicial complexes and persistent homology,” in *The 4th International Congress on Mathematical Software (ICMS)*, Hanyang University, Seoul, Korea, Aug. 2014. DOI: [10.1007/978-3-662-44199-2\\_28](https://doi.org/10.1007/978-3-662-44199-2_28).
- [112] F. Meyer, R. Overbeek, and A. Rodriguez, “Figfams: Yet another set of protein families,” *Nucleic Acids Research*, vol. 37, no. 20, pp. 6643–6654, 2009. DOI: [10.1093/nar/gkp698](https://doi.org/10.1093/nar/gkp698).
- [113] Y. Mileyko, S. Mukherjee, and J. Harer, “Probability measures on the space of persistence diagrams,” *Inverse Problems*, vol. 27, no. 12, p. 124007, 2011. DOI: [10.1088/0266-5611/27/12/124007](https://doi.org/10.1088/0266-5611/27/12/124007).

- [114] D. Morozov, *Dionysus: A C++ library for computing persistent homology*, 2012. [Online]. Available: <http://www.mrzv.org/software/dionysus/index.html>.
- [115] D. Müllner and A. Babu, *Python mapper: An open-source toolchain for data exploration, analysis and visualization*, version 0.1.13, 2013. [Online]. Available: <http://danifold.net/mapper>.
- [116] M. W. Nachman and S. L. Crowell, “Estimate of the mutation rate per nucleotide in humans.,” *Genetics*, vol. 156, no. 1, pp. 297–304, Sep. 2000.
- [117] V. Nanda, *Perseus: The persistent homology software*, version 4.0, 2015. [Online]. Available: <http://www.sas.upenn.edu/~vnanda/perseus>.
- [118] M. I. Nelson and E. C. Holmes, “The evolution of epidemic influenza,” *Nature Reviews Genetics*, vol. 8, no. 3, pp. 196–205, Jan. 2007. DOI: [10.1038/nrg2053](https://doi.org/10.1038/nrg2053).
- [119] M. I. Nelson, L. Simonsen, C. Viboud, *et al.*, “Stochastic processes are key determinants of short-term evolution in influenza A virus,” *PLoS Pathogens*, vol. 2, no. 12, e125, Dec. 2006. DOI: [10.1371/journal.ppat.0020125](https://doi.org/10.1371/journal.ppat.0020125).
- [120] H. C. Neu, “The crisis in antibiotic resistance,” *Science*, vol. 257, no. 5073, pp. 1064–1073, 1992. DOI: [10.1126/science.257.5073.1064](https://doi.org/10.1126/science.257.5073.1064).
- [121] M. E. Newman, “Modularity and community structure in networks,” *Proceedings of the National Academy of Sciences*, vol. 103, no. 23, pp. 8577–8582, Jun. 2006. DOI: [10.1073/pnas.0601602103](https://doi.org/10.1073/pnas.0601602103).
- [122] M. Nicolau, A. J. Levine, and G. Carlsson, “Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival,” *Proceedings of the National Academy of Sciences*, vol. 108, no. 17, pp. 7265–7270, 2011. DOI: [10.1073/pnas.1102826108](https://doi.org/10.1073/pnas.1102826108).
- [123] H. Ochman, J. G. Lawrence, and E. A. Groisman, “Lateral gene transfer and the nature of bacterial innovation,” *Nature*, vol. 405, no. 6784, pp. 299–304, 2000. DOI: [10.1038/35012500](https://doi.org/10.1038/35012500).
- [124] M. A. O’Malley and E. V. Koonin, “How stands the tree of life a century and a half after the origin?” *Biology Direct*, vol. 6, no. 1, pp. 1–21, 2011. DOI: [10.1186/1745-6150-6-32](https://doi.org/10.1186/1745-6150-6-32).
- [125] L. Pachter and B. Sturmfels, *Algebraic Statistics for Computational Biology*. Cambridge, UK: Cambridge University Press, Mar. 2005.

- [126] J. Penders, E. E. Stobberingh, P. H. Savelkoul, and P. F. Wolffs, “The human microbiome as a reservoir of antimicrobial resistance,” *Frontiers in microbiology*, vol. 4, 2013. doi: [10.3389/fmicb.2013.00087](https://doi.org/10.3389/fmicb.2013.00087).
- [127] C. Proux, D. van Sinderen, J. Suarez, P. Garcia, V. Ladero, G. F. Fitzgerald, F. Desiere, and H. Brüssow, “The dilemma of phage taxonomy illustrated by comparative genomics of Sfi21-like siphoviridae in lactic acid bacteria,” *Journal of bacteriology*, vol. 184, no. 21, pp. 6026–6036, Nov. 2002. doi: [10.1128/JB.184.21.6026-6036.2002](https://doi.org/10.1128/JB.184.21.6026-6036.2002).
- [128] C. Rayssiguier, D. S. Thaler, and M. Radman, “The barrier to recombination between *Escherichia coli* and *Salmonella typhimurium* is disrupted in mismatch-repair mutants,” *Nature*, vol. 342, no. 6248, pp. 396–401, Nov. 1989. doi: [10.1038/342396a0](https://doi.org/10.1038/342396a0).
- [129] J. Reininghaus, U. Bauer, S. Huber, and R. Kwitt, “A stable multi-scale kernel for topological machine learning,” in *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. arXiv: [1412.6821 \[stat.ML\]](https://arxiv.org/abs/1412.6821).
- [130] H. Rohde, J. Qin, Y. Cui, *et al.*, “Open-source genomic analysis of Shiga-toxin-producing *E. coli* O104:H4,” *New England Journal of Medicine*, vol. 365, no. 8, pp. 718–724, 2011. doi: [10.1056/nejmoa1107643](https://doi.org/10.1056/nejmoa1107643).
- [131] F. Rohwer and R. Edwards, “The phage proteomic tree: A genome-based taxonomy for phage,” *Journal of Bacteriology*, vol. 184, no. 16, pp. 4529–4535, 2002. doi: [10.1128/jb.184.16.4529-4535.2002](https://doi.org/10.1128/jb.184.16.4529-4535.2002).
- [132] F. Rohwer, M. Youle, and H. Maughan, *Life in Our Phage World*. San Diego, CA: Wholen, 2014.
- [133] N. Saitou and M. Nei, “The neighbor-joining method: A new method for reconstructing phylogenetic trees,” *Molecular Biology and Evolution*, vol. 4, no. 4, pp. 406–425, 1987. [Online]. Available: <http://mbe.oxfordjournals.org/content/4/4/406>.
- [134] G. Singh, F. Mémoli, and G. Carlsson, “Topological methods for the analysis of high dimensional data sets and 3D object recognition,” in *Eurographics Symposium on Point-Based Graphics*, The Eurographics Association, 2007, pp. 91–100.
- [135] G. J. D. Smith, D. Vijaykrishna, J. Bahl, *et al.*, “Origins and evolutionary genomics of the 2009 swine-origin H1N1 influenza A epidemic,” *Nature*, vol. 459, no. 7250, pp. 1122–1125, Jun. 2009. doi: [10.1038/nature08182](https://doi.org/10.1038/nature08182).
- [136] M. O. Sommer, G. M. Church, and G. Dantas, “The human microbiome harbors a diverse reservoir of antibiotic resistance genes,” *Virulence*, vol. 1, no. 4, pp. 299–303, 2010. doi: [10.4161/viru.1.4.12010](https://doi.org/10.4161/viru.1.4.12010).

- [137] Y. S. Song and J. Hein, “Constructing minimal ancestral recombination graphs,” *Journal of Computational Biology*, vol. 12, no. 2, pp. 147–169, 2005. DOI: [10.1089/cmb.2005.12.147](https://doi.org/10.1089/cmb.2005.12.147).
- [138] B. Sturmfels and J. Yu, “Classification of six-point metrics,” *Electronic Journal of Combinatorics*, vol. 11, R44, 2004. [Online]. Available: <http://www.combinatorics.org/ojs/index.php/eljc/article/view/v11i1r44>.
- [139] E. K. Subbarao, W. London, and B. R. Murphy, “A single amino acid in the PB2 gene of influenza A virus is a determinant of host range,” *Journal of Virology*, vol. 67, no. 4, pp. 1761–1764, Apr. 1993.
- [140] C. A. Suttle, “Marine viruses – major players in the global ecosystem,” *Nature Reviews Microbiology*, vol. 5, no. 10, pp. 801–812, Oct. 2007. DOI: [10.1038/nrmicro1750](https://doi.org/10.1038/nrmicro1750).
- [141] J. K. Taubenberger and D. M. Morens, “1918 influenza: The mother of all pandemics,” *Emerging Infectious Diseases*, vol. 12, no. 1, pp. 15–22, Jan. 2006. DOI: [10.3201/eid1209.05-0979](https://doi.org/10.3201/eid1209.05-0979).
- [142] A. Tausz, M. Vejdemo-Johansson, and H. Adams, “JavaPlex: A research software package for persistent (co)homology,” in *Proceedings of ICMS 2014*, H. Hong and C. Yap, Eds., ser. Lecture Notes in Computer Science 8592, 2014, pp. 129–136. DOI: [10.1007/978-3-662-44199-2\\_23](https://doi.org/10.1007/978-3-662-44199-2_23).
- [143] J. B. Tenenbaum, V. de Silva, and J. C. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000. DOI: [10.1126/science.290.5500.2319](https://doi.org/10.1126/science.290.5500.2319).
- [144] The Human Microbiome Project Consortium, “Structure, function and diversity of the healthy human microbiome,” *Nature*, vol. 486, no. 7402, pp. 207–214, 2012. DOI: [10.1038/nature11234](https://doi.org/10.1038/nature11234).
- [145] C. M. Thomas and K. M. Nielsen, “Mechanisms of, and barriers to, horizontal gene transfer between bacteria,” *Nature Reviews Microbiology*, vol. 3, no. 9, pp. 711–721, 2005. DOI: [10.1038/nrmicro1234](https://doi.org/10.1038/nrmicro1234).
- [146] K. Turner, Y. Mileyko, S. Mukherjee, and J. Harer. (2012). Fréchet means for distributions of persistence diagrams. arXiv: [1206.2790v2 \[math.ST\]](https://arxiv.org/abs/1206.2790v2).
- [147] L. Van der Maaten and G. E. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008. [Online]. Available: <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- [148] J. Wakeley, *Coalescent Theory*. Roberts & Company, 2009.

- [149] K. B. Walters, K. J. Emmett, J. M. Chan, D. Meroz, A. Karasin, S. Fan, G. Neumann, N. Ben-Tal, R. Rabadan, and Y. Kawaoka, “Identification of host-specific amino acids in the influenza virus PB2 polymerase subunit using machine learning approaches,” *Journal of Virology*, 2016, in preparation.
- [150] J. D. Watson and F. H. C. Crick, “Molecular structure of nucleic acids,” *Nature*, vol. 171, no. 4356, pp. 737–738, 1953. DOI: [10.1038/171737a0](https://doi.org/10.1038/171737a0).
- [151] A. R. Wattam, D. Abraham, O. Dalay, *et al.*, “PATRIC, the bacterial bioinformatics database and analysis resource,” *Nucleic Acids Research*, vol. 42, no. D1, pp. D581–D591, 2013. DOI: [10.1093/nar/gkt1099](https://doi.org/10.1093/nar/gkt1099).
- [152] B. C. Westmoreland, W. Szybalski, and H. Ris, “Mapping of deletions and substitutions in heteroduplex DNA molecules of bacteriophage lambda by electron microscopy,” *Science*, vol. 163, no. 3873, pp. 1343–1348, Mar. 1969. DOI: [10.1126/science.163.3873.1343](https://doi.org/10.1126/science.163.3873.1343).
- [153] Wikipedia, *Neighbor-joining — Wikipedia, the free encyclopedia*, 2016. [Online]. Available: [https://en.wikipedia.org/wiki/Neighbor\\_joining](https://en.wikipedia.org/wiki/Neighbor_joining).
- [154] C. R. Woese, “A new biology for a new century,” *Microbiology and Molecular Biology Reviews*, vol. 68, no. 2, pp. 173–186, Jun. 2004. DOI: [10.1128/mmbr.68.2.173-186.2004](https://doi.org/10.1128/mmbr.68.2.173-186.2004).
- [155] C. R. Woese and G. E. Fox, “Phylogenetic structure of the prokaryotic domain: The primary kingdoms,” *Proceedings of the National Academy of Sciences*, vol. 74, no. 11, pp. 5088–5090, 1977. DOI: [10.1073/pnas.74.11.5088](https://doi.org/10.1073/pnas.74.11.5088).
- [156] C. R. Woese, O. Kandler, and M. L. Wheelis, “Towards a natural system of organisms: Proposal for the domains archaea, bacteria, and eucarya.,” *Proceedings of the National Academy of Sciences*, vol. 87, no. 12, pp. 4576–4579, 1990. DOI: [10.1073/pnas.87.12.4576](https://doi.org/10.1073/pnas.87.12.4576).
- [157] World Health Organization, “Antimicrobial resistance: Global report on surveillance 2014,” Tech. Rep., 2014. [Online]. Available: <http://www.who.int/drugresistance/documents/surveillancereport/en/>.
- [158] ——, “Influenza (seasonal) fact sheet no. 211,” 2014. [Online]. Available: <http://www.who.int/mediacentre/factsheets/fs211/en/> (visited on 12/21/2015).
- [159] ——, “Cumulative number of confirmed human cases of avian influenza A(H5N1) reported to WHO,” 2016. [Online]. Available: [http://www.who.int/influenza/human\\_animal\\_interface/H5N1\\_cumulative\\_table\\_archives/en/](http://www.who.int/influenza/human_animal_interface/H5N1_cumulative_table_archives/en/) (visited on 12/21/2016).

- [160] S. Zairis, *Sakmapper: A python implementation of the mapper algorithm*, 2016. [Online]. Available: <https://github.com/szairis/sakmapper>.
- [161] S. Zairis, H. Khiabanian, A. J. Blumberg, and R. Rabadan. (2014). Moduli spaces of phylogenetic trees describing tumor evolutionary patterns. arXiv: [1410 . 0980](https://arxiv.org/abs/1410.0980) [q-bio.QM].
- [162] E. Zuckerkandl and L. Pauling, “Molecular disease, evolution, and genetic heterogeneity,” in *Horizons in Biochemistry*, M. Kasha and B. Pullman, Eds., Academic Press, 1962, pp. 189–225.
- [163] ——, “Molecules as documents of evolutionary history,” *Journal of Theoretical Biology*, vol. 8, no. 2, pp. 357–366, 1965. DOI: [10.1016/0022-5193\(65\)90083-4](https://doi.org/10.1016/0022-5193(65)90083-4).