

Outbreaks, Metastases and Homomorphisms: Phylogenetic Inference of Migration Histories of Heterogeneous Populations under Evolutionary and Structural Constraints

Kiril Kuzmin

kkuzmin1@gsu.edu

Georgia State University

Henri Schmidt

Princeton University

Sagi Snir

University of Haifa

Benjamin Raphael

Princeton University

Pavel Skums

University of Connecticut https://orcid.org/0000-0003-4007-5624

Article

Keywords: migration tree, phylogenetic inference, viral transmission, cancer metastasis

Posted Date: September 13th, 2024

DOI: https://doi.org/10.21203/rs.3.rs-5040045/v1

License: © 1 This work is licensed under a Creative Commons Attribution 4.0 International License.

Read Full License

Additional Declarations: There is **NO** Competing Interest.

Version of Record: A version of this preprint was published at Nature Communications on August 28th, 2025. See the published version at https://doi.org/10.1038/s41467-025-63411-4.

Outbreaks, Metastases and Homomorphisms: Phylogenetic Inference of Migration Histories of Heterogeneous Populations under Evolutionary and Structural Constraints

Kiril Kuzmin^{a,*}, Henri Schmidt^b, Sagi Snir^c, Ben Raphael^b, Pavel Skums^{d,*}

^aDepartment of Computer Science, Georgia State University, 25 Park Place, Atlanta, GA 30303, USA
 ^bDepartment of Computer Science, Princeton University, 35 Olden St, Princeton, NJ 08540, USA
 ^cDepartment of Evolutionary and Environmental Biology, University of Haifa, 199 Aba Khoushy Ave., Mount Carmel, Haifa 3498838, Israel
 ^dSchool of Computing, University of Connecticut, 371 Fairfield Way, Storrs, CT 06269, USA

Abstract

Many human diseases, including viral infections and cancers, are driven by the evolutionary dynamics of heterogeneous populations of genomic variants. A major type of evolutionary behavior of these populations is migration including viral transmissions and cancer metastatic spread. A common strategy for migration pathways reconstruction involves constructing a phylogenetic tree of observed genotypes and inferring its ancestral states corresponding to migration sites. Key challenges here include determining the conditions when a phylogenetic tree topology reflects the underlying migration tree structure, and balancing computational tractability, flexibility, and biological realism of inference algorithms and models.

In this study, we address these challenges using the powerful machinery of graph homomorphisms, a mathematical concept that describes how one graph can be mapped onto another while preserving its structure. We investigate how structural constraints on migration patterns and migration tree topologies influence the relationship between phylogenies and migration trees, characterize trees compatible with a given phylogeny and propose a series of algorithms to assess whether given phylogenetic and migration trees are compatible under various migration scenarios.

Leveraging our findings, we present a framework for inferring migration trees by sampling potential trees from a prior random tree distribution and identifying a subsample compatible with a given phylogeny. By varying prior tree distributions, this approach expands upon several existing models, offering a versatile strategy applicable to a variety of biological processes. We validate our methodology using simulated datasets and real data from studies of viral outbreaks and cancer metastasis, demonstrating its effectiveness across different contexts.

Keywords: migration tree, phylogenetic inference, viral transmission, cancer metastasis

1. Introduction

- 2 Many human diseases are essentially evolutionary processes.
- This includes viral infections, driven by evolving populations of
- viral variants [1], as well as cancers associated with diversify-
- ing intra-tumor subclonal lineages [2]. Although the biological
- 6 mechanisms of these diseases are different, both are fueled by
- 7 highly mutable populations of disease-causing agents, whose
- anism in RNA-dependent RNA polymerase or retroviral reverse transcriptase in RNA viruses [3, 4, 5], or from the genetic instability of tumor cells manifesting itself in somatic mutations, chromosomal gain/loss/translocation, and aneuploidy [6, 7, 8, 9]. Consequently, the general phylogenetic methodologies applied to these populations exhibit many similarities. From a

extreme genomic diversity originates from error-prone replication processes, whether due to the lack of a proofreading mech-

- methodological standpoint, they form a unique segment of phy-
- logenetics and phylodynamics, fostering a mutual exchange of

^{*}Corresponding authors.

Email addresses: kkuzmin1@gsu.edu (Kiril Kuzmin),
pavel.skums@uconn.edu (Pavel Skums)

concepts that enhances all areas of application [10, 11, 12].

A major type of evolutionary behaviour of highly mutable 53 populations is *migration*, wherein their members spread from 54 initial sites, seeding new populations at newly invaded sites. For 55 infectious diseases, such migrations equate to pathogen trans- 56 missions, whereas in cancer this process is identified with me- 57 tastatic spread. Thus, accurate inference of migration networks 58 of heterogeneous populations is crucial for public health and 59 medical research [13, 14, 15, 16].

21

22

23

24

27

31

32

33

35

The study of highly mutable population migration has been 61 significantly enhanced by groundbreaking advancements in se- 62 quencing technologies. State-of-the-art high-throughput targe- 63 ted sequencing and single-cell DNA sequencing enable the cap- 64 ture of detailed population snapshots at exceptionally high res- 65 olutions, facilitating fine-grained analysis down to the level of 66 individual genotypes [16, 17, 18, 19]. In particular, this allows 67 to examine population migration on the level of individual mi- 68 gration events [20, 17, 21, 22].

A wide array of methods has been developed specifically 70 for reconstructing viral and bacterial transmission trees, reflect-71 37 ing the substantial interest in tracking the spread of infectious 72 diseases. The arsenal of tools available for reconstructing trans-73 mission networks is extensive, including but not limited to Out- 74 breaker, Outbreaker 2 [23, 24], SeqTrack [25], SCOTTI [26], 75 SOPHIE [27], Phybreak [28], Bitrugs [29], BadTrIP [30], Phy- 76 loscanner [31], StrainHub [32], TransPhylo [33, 34] (along with 77 its extension TransPhyloMulti [35]), STraTUS [36], TreeFix-78 TP [37], QUENTIN [38], VOICE [39], HIVTrace [40], GHOST 79 [41], MicrobeTrace [42], SharpTNI [43], TiTUS [12], TNeT 80 [44], AutoNet [45], and others [46, 47, 34, 48, 49, 50, 51, 52, 81 53]. These tools have been instrumental in investigation of out- 82 breaks and monitoring the transmission dynamics of pathogens 83 like HIV, hepatitis C (HCV), SARS, MERS and SARS-CoV-2 84 [54, 55, 56, 57, 20, 58].

Similarly, the development of methods for deducing metastatic spread histories is burgeoning, driven by advancements in single-cell DNA sequencing and CRISPR-based lineage tracing technologies. Currently, the repertoire of tools in this domain includes MACHINA [22], FitchCount (as part of the Cassiopeia suite) [21, 59], PathFinder [60], TCC, PCC, and PCCH [61]. Notwithstanding the relatively shorter list, the impact of these tools is growing, with several studies published recently leveraging these methods to gain insights into the mechanisms of metastatic spread [21, 62, 63, 64].

Despite significant progress in the field, the wide variety of existing methods underscores that the challenge of accurately inferring heterogeneous population migration remains unresolved. This diversity of approaches indicates both the complexity of the problem and the ongoing efforts to refine and improve upon existing methods. Additionally, a major barrier to advancement in the field is the relative isolation of viral and cancer genomics fields. Some of the aforementioned tools are based on conceptually similar techniques - this applies, for example, to STraTUS (which focus on viral transmission) and FitchCount (which addresses metastatic spread) that, as demonstrated in this paper, yield virtually identical results when applied to the same datasets. This lack of interdisciplinary exchange of ideas often leads researchers to inadvertently duplicate efforts, thereby impeding progress in both fields.

Phylogenetics and phylodynamics provide the most widely used methodological frameworks for migration tree/network reconstruction [35, 22]. However, their application in this context is not straightforward. A phylogenetic tree does not directly equate to a migration tree [35], as the nodes in a phylogenetic tree represent divergences of lineages rather than specific migration events [35, 65]. While some of these divergences may result from migration, others occur within previously invaded sites. Therefore, deriving a migration tree from a phyloge-

netic tree essentially requires solving an ancestral trait infer-120 ence problem, wherein the internal nodes of a phylogenetic tree121 are annotated with labels that indicate whether each divergence event occurred within a site or as a result of seeding of a new site upon migration. Furthermore, it is crucial to effectively leverage the full spectrum of intra-site (within-host or within-125 tumor) population diversity uncovered by high-throughput sequencing, that often provide a strong signal for migration in-127 ference. For example, the paraphyletic relationships between populations suggest recent migrations between corresponding sites [66, 67, 39, 68].

The challenges mentioned above highlight several crucial 13 questions that remain unresolved and warrant further exploration. These questions include:

100

101

102

103

104

106

107

108

109

110

111

112

113

114

117

118

119

1) To what extent and under which conditions does the topol-134 ogy of a phylogenetic tree reflect the structure of the underly-135 ing migration tree? This question becomes particularly impor-136 tant when the level of paraphyly in the labeled phylogeny is 137 low, a situation not uncommon as the paraphyletic signal tends₁₃₈ to diminish over time and with smaller sample sizes [66]. A₁₃₉ number of studies have looked at this subject but their conclu-140 sions were mixed. Some studies have found that certain mi-141 gration patterns, such as super-spreading, migration chains, or,142 more generally, migrations within networks formed by differ-143 ent models like Erdös-Rényi [69], preferential attachment [70],144 or Watts-Strogatz models [71], lead to quantitatively distinct₁₄₅ phylogenetic tree topologies [72, 73, 74]. Additionally, the 146 spatial structure and dynamics of heterogeneous populations, 147 which are directly related to migrations pathways, have been 148 shown to affect the phylogeny structure of both viral and tumor₁₄₉ populations [75, 76, 77, 78]. On the other hand, other studies₁₅₀ report that the direct impact of migration patterns on phyloge-151 netic trees ranges from minimal to moderate [79, 80, 81, 82].152 Finally, certain studies have drawn mixed conclusions, indicat-153 ing that while some migration characteristics are reflected in the phylogeny, others are not [83].

2) How can we limit the solution space to balance computational feasibility, accuracy of inference, generalizability, and biological realism? The space of migration trees compatible with given phylogenetic trees is often vast, and its properties are not well understood [36, 12]. A sampling-consensus approach is one method to address solution ambiguity, where feasible migration trees are sampled and summarized in a weighted consensus graph, with weights reflecting posterior probabilities of edges [44, 12, 43, 26, 21]. However, the size of solution space may restrict the depth of sampling. As a response, it is common practice to narrow down the solution space to a set of plausible migration trees optimizing a specific objective function under evolutionary-based constraints. Employing constrained models also aids in preventing overfitting in presence of missing data and errors.

Various objectives and constraints have been implemented by existing methods. Limiting number of migration events, sizes of bottlenecks or numbers of back-migrations [31, 12, 44, 22, 59, 21, 61, 36, 25] is more computationally efficient and scalable due to utilization of dynamic programming [84, 85], making such approaches practical in both molecular epidemiology and computational oncology. These can also be formulated as Integer Linear Programming (ILP) problems [86] and solved with reasonable efficiency using existing ILP solvers. Models with more complex Bayesian objectives with constraints regularized as priors [35, 28, 26, 23, 30] offer a richer, biologically nuanced perspective but suffer from scalability issues and usually rely on generic methods like Markov Chain Monte Carlo (MCMC) sampling, which may not yield optimal solutions, in part due to a lack of problem-specific mathematical strategies. Balancing computational efficiency with biological comprehensiveness presents a notable challenge, compounded by the uncertainty of how much constraints and objectives truly limit the 188 solution space [36, 12].

3) How to incorporate a variety of models for phylogenetic ¹⁹⁰ inference of migrations into a unified modular computational ¹⁹¹ framework?

156

157

158

159

160

161

166

167

168

169

170

171

172

174

175

176

177

178

179

180

181

182

185

186

187

Migration inference models draw on varied biological or 193 epidemiological assumptions. For instance, viral transmission 194 inference often incorporates case-specific temporal data like infectious periods, exposure intervals, symptom onset, diagnosis 196 or sample collection dates to establish order of infections and 197 eliminate unlikely transmission links [23, 47, 28, 33, 12, 26, 198] 29]. In some rare cases, contact networks are known and can be 199 used for the same purpose [87, 58]. While effective, such data²⁰⁰ is often unavailable, non-informative, or sensitive, particularly 201 for endemic and pandemic diseases caused by HIV, Hepatitis²⁰² C, SARS-CoV-2, or Influenza [27, 38, 53]. In situations where 203 case-specific data cannot be used, genomic epidemiology tools²⁰⁴ resort to broader assumptions, like the expected degree distri-205 bution of transmission networks implied by a structure of a sus-206 ceptible population [38, 27]. Similarly, methods for inferring²⁰⁷ metastatic spread use constraints defined by so-called migration 208 patterns that reflect realistic cancer migration scenarios, such²⁰⁹ as monoclonal, polyclonal or multi-source seeding [22, 61]. A²¹⁰ commonality across all these methods is the use of structural²¹¹ constraints on feasible transmission networks that are consid-212 ered as subgraphs of a larger "pattern" graph. It suggests a²¹³ need for a versatile, modular migration inference framework 214 that integrates these varied approaches on a unified algorithmic 215 and mathematical basis, akin e.g. to the BEAST framework for 216 Bayesian evolutionary analysis [88].

Addressing these challenges requires a comprehensive in-²¹⁸ vestigation into the mathematical properties of migration trees²¹⁹ compatible with a given phylogeny, a topic that is not yet fully understood [36]. Our study aims to advance this area by devel-²²¹

oping a novel methodology based on powerful techniques from graph theory and combinatorics. Additionally, we will use this methodology to introduce novel algorithmic approaches for inferring migration trees.

A number of earlier studies has achieved important progress in this area. Several studies noticed that migration trees compatible with a given phylogeny correspond to partitions of the phylogeny's node set or to coloring of its branches [51, 36, 33, 34]. These observations have informed the development of methods to enumerate and sample these trees, assuming a complete migration bottleneck [36] and a known sequence of migrations [89]. In fact, as we argue in this paper, the relation between a phylogenetic tree and a migration tree is described by the concept of a *graph homomorphism* that generalize both partitions and colorings.

Graph homomorphism is essentially a mapping between the vertices of two graphs that preserve their structure [90]. The theory of graph homomorphisms is well-established area of discrete mathematics, with deep results and rich methodology. All types of migration trees discussed so far can be described by a graph homomorphism with specific constraints. For example, migration trees compatible with a given phylogeny under the assumptions of complete sampling and complete bottleneck, that has been studied in [51, 36] are *minors* [91, 92] of the phylogeny or, equivalently, its homomorphic images such that inverse images of all vertices are connected subtrees.

We use mathematical and algorithmic machinery of graph homomorphism theory to delve into the details of migration inference. Specifically, we provide necessary and sufficient conditions describing trees compatible with a given phylogeny and propose a series of algorithms that evaluate the compatibility of phylogenetic and migration trees under various evolutionary scenarios through the construction of corresponding homomorphisms. We examine particular structural constraints on migration patterns and migration tree topologies to understand their₂₅₅ influence on the relationship between a phylogeny and a migra-₂₅₆ tion tree.

Based on aforementioned insights, we propose a general₂₅₈ framework for migration inference that samples candidate mi-259 gration trees from a chosen prior random tree distribution, and₂₆₀ identifies a subsample of trees compatible with a given phy-261 logeny. By varying prior tree distributions, this approach ex-262 pands upon and generalizes several existing models, offering a₂₆₃ versatile and computationally efficient strategy applicable to a264 variety of biological processes associated with heterogeneous265 population migration. Crucially, the proposed framework is266 computationally fast, enabling biomedical and public health re-267 searchers to quickly test different tree priors that represent com-268 mon migration models. Beyond migration reconstruction, pro-269 posed methods can be used for investigating how phylogenies₂₇₀ constrain the space of possible migration trees, for inferring an₂₇₁ order of known or suspected migration events, and for deter-272 mining potential migration events that are definitively ruled out₂₇₃ by a phylogeny [36].

Proposed methodology was validated using both simulated₂₇₅ datasets and real experimental data gathered from studies of vi-₂₇₆ ral outbreaks and cancer metastasis, demonstrating its effective-₂₇₇ ness and applicability across different contexts.

2. Methods

225

226

227

228

229

230

231

232

236

237

238

239

240

241

249

250

251

252

253

2.1. Basic definitions

Throughout this paper, we consider a pair of trees: a phylogenetic tree and a migration tree. For clarity, we refer to elements of a phylogeny as *nodes*, and to elements of a migration network as *vertices*, and denote them by Greek and Latin letters, respectively.

The problem of *migration network inference* is set up as follows. The input is a phylogenetic tree $\Psi = (V(\Psi), E(\Psi))$,

with the leaf set $L(\Psi)$ representing genomic variants belonging to different subpopulations (or *demes*), denoted by \mathcal{L} . The tree Ψ can be a standard binary phylogeny or non-binary mutation tree used in most cancer studies. Each leaf $\lambda \in L(\Psi)$ has an assigned site label (or color) $l_{\lambda} \in \mathcal{L}$. The aim is to expand this labeling from the leaves to all nodes in the tree, creating a full labeling $f: V(\Psi) \to \mathcal{L}$. In this model, any multi-colored tree edge $\alpha\beta$ represents a migration of genomic variants between demes $f(\alpha)$ and $f(\beta)$. The migration tree $T = T(\Psi, I)$ and with vertices $V(T) = \mathcal{L}$, is then formed by contracting the nodes with the same color [51].

As mentioned in the introduction, researchers often seek migration trees satisfying particular constraints restricting types of migration or tree topologies. These constraints can be encoded using a *transition pattern graph G* that describes permissible patterns of migration (specific examples are provided in the following subsection). We will first consider the situation when G is a simple graph; later on, it will be extended to the cases when G is a random graph characterized by some probability distribution. In this model, a migration tree should be isomorphic (i.e. identical up to relabeling of vertices) to a subgraph of the transition pattern. Any corresponding labeling will be called *feasible*.

The relations between the phylogeny Ψ , the migration tree T and the transition pattern G can be captured using the concept of a *graph homomorphism*. A homomorphism $f: \Psi \to G$ [90] is an adjacency-preserving mapping between vertex sets of these graphs, i.e. $f(u)f(v) \in E(G)$ if $uv \in E(\Psi)$. For the sake of mathematical rigor, here and throughout this paper we assume that a transition pattern is *reflexive*, i.e. every vertex is adjacent to itself. With this condition in place, any feasible labeling f is a homomorphism from Ψ to G, making the migration inference problem essentially a problem of finding such a homomorphism. In graph theory, this type of problems is sometimes

279

280

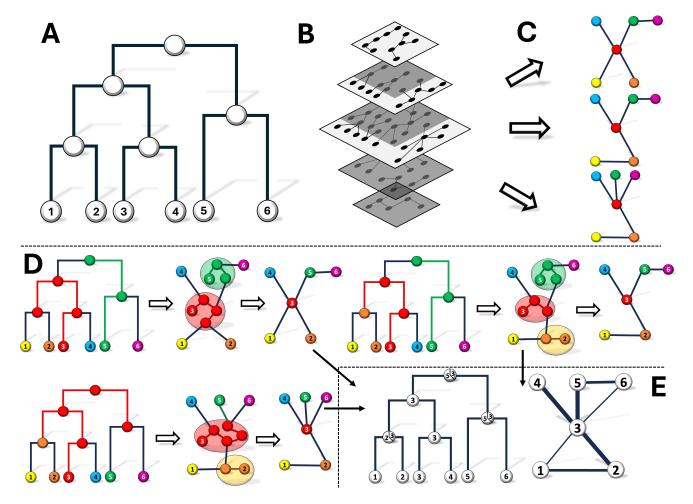


Figure 1. SMiTH: Sampling MIgration Trees with Homomorphisms. A: Input phylogenetic tree. B: Distribution of possible migration trees. Parallel rectangles depict a probability density function, with each rectangle's width proportional to the corresponding probability. In practice, the distribution is represented either by a random graph model or by a stochastic graph generation procedure. C: Candidate migration trees sampled from the distribution B. D: Homomorphisms from the phylogeny to three sampled trees. In the phylogenies, nodes are color-coded by their homomorphic images in a migration tree. The phylogeny layouts in the middle of each subfigure showcases how homomorphism transforms them into sampled trees. E: consensus solution derived from homomorphisms in D. The solution is shown as potential color distributions for the phylogeny's nodes (left) or as a graph where possible migration edges are weighted according to the number of supporting solutions (right), with the edge thickness indicating weight.

referred to as an *G-coloring* of the tree Ψ [90].

289

Throughout this paper, we use standard graph theory no-301 tations. We denote by Ψ_{α} a subtree of Ψ rooted at the node302 α . For any graph G, G[X] represents the subgraph induced by303 the subset of vertices X. We use $u \sim v$ to indicate that ver-304 tices u and v are adjacent. The set of neighbors of a vertex u305 in G is denoted as $N_G(u)$, $\deg_G(u) = |N_G(u)|$ is the degree of 306 u, $N_G[u] = N_G(u) \cup \{u\}$ and $N_G[X]$ is the union of the neigh-307 borhoods of all vertices in a set X. Additionally, the distance between any two vertices u and v in G is denoted by $d_G(u, v)$,

and $P_G(u, v)$ refers to the corresponding shortest path between

them. When the graph is clear from a context, we may omit the subscripts in these notations. In the case of a phylogeny Ψ , all distances and paths are *undirected*.

Additionally, to simplify the notation, we will apply settheoretical operations (e.g. intersection and union) directly to subgraph of G, with the understanding that the resulting subgraph is induced by the set obtained from applying these operations to the vertex sets of the original subgraphs.

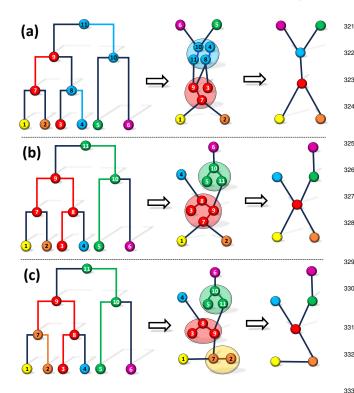


Figure 2. Solutions for the Migration History Inference Problem under various constraints. For each scenario, a phylogenetic tree is illustrated in two different layouts on the left and in the middle, with the corresponding migration 334 tree displayed on the right. In the phylogenies, nodes are color-coded by their homomorphic images in the migration tree. The layout in the middle showcases how homomorphism transforms a phylogeny into a migration tree. (a) Uncon-335 strained solution. Subtrees formed by blue and red nodes are not connected, indicating a violation of the convexity constraint. (b) Convex solution. Each336 color-coded subtree is connected, but compactness is violated in the subtree rooted at node 7. (c) Convex and compact solution.

2.2. Migration inference under structural constraints

313

314

315

316

317

318

319

320

The simplest form of the migration inference challenge appears when the vertices of the transition pattern graph G directly correspond to the labels of the phylogenetic tree leaves, that is, $V(G) = \mathcal{L}$. This variant is referred to as *labeled inference*. The transition pattern G can be an unweighted graph or a random graph with specified edge probabilities. Several scenarios exemplify this problem:

- For viral transmissions, *G* could reflect a contact network of potential hosts [87, 58], where an infection spreads through direct interactions within this network.
- Another viral transmission model assumes that transmis-350
 sion is only possible between hosts with overlapping ex-351

posure intervals [30, 12]. In this case *G* is an *interval* graph [93], i.e. a graph with vertices representing time intervals and edges connecting vertices whose intervals intersect.

 For metastatic spread inference, G may represent a circulatory network [94, 95, 96], with vertices representing organs and edges reflecting the prior probabilities of cancer spreading between organs.

Problem 1 (labeled migration inference)

Input:

- (a1) a phylogenetic tree Ψ with leaf labels $(l_{\lambda})_{\lambda \in L(\Psi)}$ forming the label set \mathcal{L} ;
- (b1) a transition pattern G with $V(G) = \mathcal{L}$;

Output:

337

338

(c1) a homomorphism $f: \Psi \to G$ such that $f(\lambda) = l_{\lambda}$ for every $\lambda \in L(\Psi)$.

In practical settings, however, Problem 1 might be too restrictive or not fully reflective of reality. The main issue is the absence of a known mapping between the leaf labels and the vertices of the transition pattern. In other words, the pattern G represents the permissible topology of migration rather than specific allowed migration links. For instance, we may expect that a viral transmission network likely includes a superspreader, but the exact subpopulation associated with this superspreader is not known. The following models to this variant of the problem:

Viral transmission: The transition pattern G could represent a random graph that describes expected characteristics of transmission trees, like being scale-free or having a particular expected degree distribution. Such models draw on known expected properties of contact networks

essential for infection spread [97, 98, 99, 100] without₃₈₃ requiring actual host contact information, and have been₃₈₄ explored in several studies [38, 27].

• Metastatic spread inference: In this case, a transition pattern *G* may describe plausible evolutionary scenarios of cancer migration, such as monoclonal or polyclonal seding from single or multiple sources [22, 61]. A related idea has been applied in studies analyzing CRISPR-390 based lineage tracing phylogenies, where *G* specifies a so-called star homoplasy model [101].

In this version of migration inference problem, is no ex-393 362 plicitly given mapping between the set of leaf labels and the394 363 vertices of the transition pattern. Instead, a feasible homomor-395 364 phism should map leaves with distinct labels to distinct ver-396 365 tices of G and vice versa, i.e. for any $\lambda_1, \lambda_2 \in L(\Psi)$ we have₃₉₇ 366 $f(\lambda_1) \neq f(\lambda_2)$ whenever $l_{\lambda_1} \neq l_{\lambda_2}$ (Fig. 2 (a)). We will call a₃₉₈ homomorphisms satisfying this requirement a label-distinctive399 homomorphism. It should be noted that finding such a homo-400 morphism seems to be a non-standard variant of a graph homo-401 morphism problem, that, to the best of our knowledge, has not₄₀₂ been studied previously. 372

In such context, the first question to be asked is whether a_{404} given subtree T of the transition pattern G is compatible with a_{405} the phylogeny Ψ , i.e. whether there exist a label-distinctive ho- a_{406} momorphism $\Psi \to T$.

Problem 2 (unlabeled migration inference)

378 *Input*:

373

374

375

376

353

358

359

360

361

- (a2) a phylogenetic tree Ψ with leaf labels $(l_{\lambda})_{\lambda \in L(\Psi)}$;
- (b2) a tree T;
- 381 Output:
- (c2) a label-distinctive homomorphism $f: \Psi \to T$.

A more restricted version of Problem 3 includes an additional *convexity constraint* [102, 103], where tree nodes mapping to the same vertex of G form a connected subtree of Ψ (Fig. 2 (b)):

Problem 3 (convex unlabeled migration inference)

Input: (a2) and (b2)

Output:

(c3) a label-distinctive homomorphism $f: \Psi \to T$ such that induced subgraphs $\Psi[f^{-1}(v)]$ are connected for all $v \in V(T)$.

We will refer to a homomorphism satisfying (c3) as a *convex homomorphism*. Convex homomorphisms to trees describe migrations with a complete bottleneck; such migration trees were the focus of extensive research in previous studies [51, 36, 89].

To define another type of constraints, consider a labeling l of nodes of the phylogeny Ψ . The labeling l is *compact* if the label of the most recent common ancestor (MRCA) for any group of leaves matches one of the labels within that group.

The rationale for this is based on the understanding that migrations within any given subtree could follow one of the following scenarios: either (i) migrations involve only the demes represented by the leaves of the subtree, or (ii) migrations include external demes, but lineages from these demes were not sampled due to extinction or incomplete collection. Although both scenarios are feasible, the first is more parsimonious, especially when sampling is dense and migration events span a short period of time. Hence, homomorphisms that align with this compact labeling are referred to as compact.

Problem 4 (compact unlabeled migration inference)

- 2 Input: (a2) and (b2)
- Output:
 - (c4) a label-distinctive homomorphism $f: \Psi \to T$ such that for any node $\alpha \in V(\Psi)$ we have $f(\alpha) \in f(L(\Psi_{\alpha}))$.

414

446

Finally, it is often desirable to produce a sample of potential solutions rather than a single solution. Such approach offers statistical backing for potential migration network edges and 418 logically shifts the migration inference problem to a Bayesian 419 paradigm. Potential solutions can be sampled from a random 420 graph distribution or, if a transition pattern is a deterministic 421 graph, from the uniform distribution of its subgraphs with spec-422 ified properties (e.g. subtrees). The sampling-based version of 423 the migration inference problem can be formulated as follows: 424

25 Problem 5 (unlabeled migration sampling)

426 *Input*:

- (a5) a phylogenetic tree Ψ with leaf labels $(l_{\lambda})_{\lambda \in L(\Psi)}$;
- (b5) a random transition pattern G with edge probabilities p: $E(G) \to [0,1] \text{ that define a distribution } \mathcal{D} \text{ of subtrees of}$

431 *Output*:

436

432 (c5) a sample $S = \{T_1, \dots, T_m\}$ from the distribution \mathcal{D} , whe433 re each subtree T_i is a homeomorphic image of Ψ , and
434 the corresponding label-distinctive homomorphisms f_i :
445 $\Psi \to T_i$ are possibly convex and/or compact.

2.3. Labeled migration inference

Problem 1 can be efficiently solved in polynomial time using
dynamic programming, as detailed in Algorithm 1. Despite its452
simplicity, we include the algorithm here due to its relevance453
for subsequent, more complex methods.

454

2.4. Unconstrained unlabeled migration inference

A possible strategy to solve Problem 2 involves identifying a bijection $g: \mathcal{L} \to V(T)$ that can be extended to a homomorphism $\Psi \to T$ using Algorithm 1. We will describe such bijections as
feasible.

Algorithm 1 Unlabeled migration inference

```
1: Let \rho be the root of \Psi.
                                      Perform a post-order
    traversal of the tree \Psi.
    for every node \alpha of the traversal do
 3:
         construct a set of potential images I(\alpha) \subseteq \mathcal{L}:
 4:
         if \alpha is a leaf then
 5:
             I(\alpha) \leftarrow \{l(\alpha)\};
 6:
         else
            Suppose that \beta_1, \ldots, \beta_k are children of \alpha.
 7:
    Then I(\alpha) consists of vertices x \in V(G) such that
    there exist vertices y_i \in I(\beta_i), i = 1, k such that
    y_i \sim x.
 8:
        end if
9:
    end for
10:
     if I(\rho) = \emptyset then
11:
         homomorphism f does not exist
12:
     else
         perform a pre-order traversal of \boldsymbol{\Psi} and
13:
     construct a homomorphism f as follows
         for every node \alpha of the pre-order do
14:
             if \alpha = \rho then
15:
                 select any v \in I(\rho) and set f(\rho) \leftarrow v;
16:
17:
                 Let \omega be the parent of \alpha. Choose v \in I(\alpha)
     such that v \sim f(\omega) and set f(\alpha) \leftarrow v.
19:
             end if
20:
         end for
21: end if
```

Let $L_g(u) = \{\lambda \in L(\Psi) : g(l_{\lambda}) = u\}$ be the set of leaves in Ψ whose labels map to u. The following theorem establishes a necessary and sufficient condition for the bijection feasibility.

Theorem 1. The bijection $g: \mathcal{L} \to V(T)$ is feasible if and only if

$$d_T(u_1, u_2) \le d_{\Psi}(\lambda_1, \lambda_2) \tag{1}$$

for all $u_1, u_2 \in V(T)$, $\lambda_1 \in L_g(u_1)$, $\lambda_2 \in L_g(u_2)$.

Proof. The necessity of the condition stated in the theorem is a known property of graph homomorphisms [90]. However, it is generally not sufficient [90]. For our specific type of the homomorphism problem, we will demonstrate that it indeed suffices.

Consider a bijection g satisfying the condition (1). Define $C_{\alpha} \subseteq V(T)$ as the image set of its clade, that is, $C_{\alpha} = g(l_{\lambda}) : \lambda \in L(\Psi_{\alpha})$. Additionally, for a node $\alpha \in V(\Psi)$ and a vertex $u \in V(T)$, define $B_{\alpha}(u)$ as a ball centered at u in T, with

493

radius

$$r_{\alpha,u} = \min_{\lambda \in L_{\alpha}(u) \cap L(\Psi_{\alpha})} d_{\Psi}(\alpha, \lambda)$$
 (2)

i.e.,
$$B_{\alpha}(u) = \{ v \in V(T) : d_T(u, v) \le r_{\alpha, u} \}.$$

We proceed by establishing two auxiliary facts. The first follows directly from the properties of a tree:

Lemma 1. Let $X_1, ..., X_k$ be subsets of vertices of the tree₄₈₆

T such that each subsets induces a connected subgraph and₄₈₇ $\bigcap_{i=1}^k X_i \neq \emptyset$. Then $N[\bigcap_{i=1}^k X_i)] = \bigcap_{i=1}^k N[X_i]$.

Suppose now that $I(\alpha): \alpha \in V(\Psi)$ are the sets of potential node images produced by Algorithm 1, given the matching of leaf labels in Ψ and vertices of T through the bijection g. These sets are described by the following lemma.

Lemma 2. For every node $\alpha \in V(\Psi)$, $I(\alpha) = \bigcap_{u \in C_{\alpha}} B_{\alpha}(u)$.

Proof. The proof of the lemma proceeds by induction. Consider a node $\alpha \in V(\Psi)$. The lemma's assertion is trivially true when α is a leaf. Assume now that α is an internal node with children β_1, \ldots, β_k , each set $I(\beta_i)$ is non-empty and, by the inductive assumption,

$$I(\beta_i) = \bigcap_{u \in C_{\beta_i}} B_{\beta_i}(u). \tag{3}$$

Then $C_{\alpha} = \bigcup_{i=1}^{k} C_{\beta_i}$; furthermore, according to Algorithm 1⁵⁰² and the equality (3) we have

$$I(\alpha) = \bigcap_{i=1}^{k} N[I(\beta_i)] = \bigcap_{i=1}^{k} N\left[\bigcap_{u \in C_{\beta_i}} B_{\beta_i}(u)\right]. \tag{4}$$

Now consider a vertex $u \in C_{\alpha}$. Let $\beta(u)$ be the child of α^{504} such that $r_{\beta(u),u} = \min_{i=1}^k r_{\beta_i,u}$; in cases where there are multiple⁵⁰⁵ such children, we pick any of them. Consequently, we have⁵⁰⁶ $r_{\alpha,u} = r_{\beta(u),u} + 1$, leading to the relation $B_{\alpha}(u) = N[B_{\beta(u)}(u)]$.⁵⁰⁷ Furthermore, it is obvious that $B_{\beta(u)}(u) \subseteq B_{\beta_i}(u)$ for all nodes⁵⁰⁸ β_i such that $u \in B_{\beta_i}(u)$. Together, these observations imply the 509

following sequence of equalities:

$$\bigcap_{u \in C_{\alpha}} B_{\alpha}(u) = \bigcap_{u \in C_{\alpha}} N[B_{\beta(u)}(u)] = \bigcap_{i=1}^{k} \bigcap_{u : \beta(u) = \beta_{i}} N[B_{\beta_{i}}(u)] =$$

$$= \bigcap_{i=1}^{k} \bigcap_{u \in C_{\beta_{i}}} N[B_{\beta_{i}}(u)] \quad (5)$$

By Lemma 1, $N\left[\bigcap_{u\in C_{\beta_i}} B_{\beta_i}(u)\right] = \bigcap_{u\in C_{\beta_i}} N[B_{\beta_i}(u)]$, and thus the expressions (4) and (5) are equal. This completes the proof of Lemma 2.

According to Lemma 2, Algorithm 1 succeeds whenever $\bigcap_{u \in C_{\alpha}} B_{\alpha}(u) \neq \emptyset$ for every node $\alpha \in V(\Psi)$. To establish that this condition holds, we invoke so-called *Helly property* of subtrees. A family of sets S_1, \ldots, S_k has a Helly property [104] if $\bigcap_{i=1}^k S_i \neq \emptyset$ whenever $S_i \cap S_j \neq \emptyset$ for every $i, j \in [k]$; in other words, the existence of non-empty pairwise intersections guarantees a non-empty total intersection.

Subtrees of a given tree are known to have the Helly property [105]. The subsets $B_{\alpha}(u)$ obviously induce subtrees of T. Therefore, to prove the theorem, it is sufficient to demonstrate that for every node $\alpha \in V(\Psi)$ and for every pair of vertices $u_1, u_2 \in C_{\alpha}$, the intersection $B_{\alpha}(u_1) \cap B_{\alpha}(u_2)$ is non-empty.

Select two leafs $\lambda_1, \lambda_2 \in L(\Psi_\alpha)$ that minimize the distance between nodes of the sets $L_g(u_1) \cap L(\Psi_\alpha)$ and $L_g(u_2) \cap L(\Psi_\alpha)$. According to the theorem's conditions, we have:

$$d_T(u_1, u_2) \le d_{\Psi}(\lambda_1, \lambda_2) \le d_{\Psi}(\alpha, \lambda_1) + d_{\Psi}(\alpha, \lambda_2). \tag{6}$$

This implies the existence of a path between u_1 and u_2 in T with a length at most $d_{\Psi}(\alpha, \lambda_1) + d_{\Psi}(\alpha, \lambda_2)$. On this path, there is at least one vertex v such that $d_T(u_1, v) \leq d_{\Psi}(\alpha, \lambda_1)$ and $d_T(u_2, v) \leq d_{\Psi}(\alpha, \lambda_2)$. Consequently, v belongs to both $B_{\alpha}(u_1)$ and $B_{\alpha}(u_2)$, thereby completing the proof.

Theorem 1 establishes that the Unlabeled Migration Inference problem is algorithmically equivalent to the problem of finding a bijection that satisfies (1). First of all, this allows us535 to demonstrate that Problem 2 is \mathcal{NP} -hard. We will prove it₅₃₆ through a reduction from the following problem:

Graph Bandwidth problem.

Given: A graph G. 515

Find: The minimal integer K = bw(G) for which there exists is 539 516

a bijection $f: V(G) \to \{1, \dots, |V(G)|\}$ such that 517

$$|f(u) - f(v)| \le K \text{ for all } u \sim v.$$
 (7)⁵⁴²

The Graph Bandwidth problem is \mathcal{NP} -hard [106], and, mo-518 reover, it cannot be approximated within any constant factor 519 unless $\mathcal{P} = \mathcal{NP}$, even when the input graph G is a tree [107]. 520 **Theorem 2.** The Unlabeled Migration Inference problem is₅₄₇ \mathcal{NP} -hard, even when all leaf labels in the phylogeny Ψ are ₅₄₈ unique.

Proof. For our purposes, it is more convenient to use an equivalent condition for Graph Bandwidth problem:

$$|f(u) - f(v)| \le Kd_G(u, v) \text{ for all } u, v \in V(G).$$
 (8)⁵⁵⁰

The fact that (8) implies (7) is obvious. To demonstrate⁵⁵² that (7) implies (8), consider the shortest (u, v)-path in G(u) $x_0, x_1, ..., x_{d-1}, x_d = v$), where $d = d_G(u, v)$. Then we have

$$|f(u) - f(v)| = |\sum_{i=1}^{d} (f(x_{i-1}) - f(x_i))| \le Kd.$$

$$\le \sum_{i=1}^{d} |(f(x_{i-1}) - f(x_i))| \le Kd.$$
 This implies that $bw(S) \le \frac{5}{3}K^*$.

Now suppose that T' is an input tree of Graph Bandwidth 556 529 problem. We can assume that $bw(T') \ge 3$, since graphs where start $bw(T') \ge 3$, since graphs where $bw(T') \le 2$ are recognizable in linear time [108]. To construct an instance for Problem 2, for an integer $K \ge 3$, we proceed as follows:

1) The input phylogeny Ψ_K is constructed by (a) subdivid-

534

ing every edge of T' into K edges; (b) attaching a leaf labeled u' to every node $u \in V(T')$.

2) The input tree T is an n-vertex path P_n with V(T) = $\{1, \ldots, n\}$

For the trees constructed in this manner we have:

(a)
$$L(\Psi_K) = \{u' : u \in V(T')\};$$

(b)
$$d_{\Psi_K}(u', v') = Kd_S(u, v) + 2;$$

(c)
$$d_T(i, j) = |i - j|$$
.

We will demonstrate that a polynomial-time algorithm for Problem 2 leads to a $\frac{5}{3}$ -approximation algorithm for the Graph Bandwidth problem. Assume the existence of such an algorithm for Problem 2. Let K^* be the smallest integer for which this algorithm produces a sought-for homomorphism $\Psi_{K^*} \to T$. To establish the $\frac{5}{3}$ approximation factor, we need to demonstrate the following relationship:

$$K^* \le bw(T') \le \frac{5}{3}K^* \tag{9}$$

To establish an upper bound, let us consider the feasible bijection $g^*: L(\Psi_{K^*}) \to \{1, \dots, n\}$. We extend this bijection to the nodes of V(T') by setting $g^*(u) = g^*(u')$. Given the inequality (1) and assuming that $K^* \ge 3$, we have:

$$\begin{split} |g^*(u) - g^*(v)| &= d_T(g^*(u'), g^*(v')) \le d_{\Psi_{K^*}}(u', v') = \\ &= K^* \cdot d_{T'}(u, v) + 2 \le \frac{5}{3} K^* d_{T'}(u, v). \end{split}$$

Conversely, if $bw(T') = K < K^*$, and the mapping g: $V(T') \rightarrow \{1, \dots, n\}$ is the bijection reflecting this bandwidth, then we can extend it to $L(\Psi_K)$ by setting g(u') = g(u). This yields:

586

599

601

$$\begin{split} d_T(g(u'),g(v')) &= |g(u)-g(v)| \leq K d_{T'}(u,v) < \\ &< K d_{T'}(u,v) + 2 = d_{\Psi_K}(u',v'). \end{split}$$

This inequality contradicts the assumption that K^* is minimal. Therefore, we must have $bw(S) \ge K^*$. This completes the proof.

Although the Unlabeled Migration Inference problem is \mathcal{NP} hard, we can use Theorem 1 to approach it using Integer Linear₅₈₇

Programming (ILP). We define a feasible bijection $g: \mathcal{L} \rightarrow_{588}$ V(T) using binary variables x_{iu} , where $x_{i,u} = 1$ if g(i) = u, for g(i) = u and g(i) = u for g(i) = u

$$\sum_{i \in \mathcal{I}} \sum_{u \in V(T)} \delta(i) \deg(u) x_{iu} \to \max \tag{10}$$

s.t.

569

$$\sum_{u \in V(T)} x_{iu} = 1, \quad i \in \mathcal{L}; \tag{11}$$

$$\sum_{i \in \mathcal{L}} x_{iu} = 1, \quad u \in V(T); \tag{12}^{597}$$

$$x_{iu} + x_{iv} \le 1$$
, $i, j \in \mathcal{L}, u, v \in V(T)$

and
$$\min_{\lambda_i \in l^{-1}(i), \lambda_j \in l^{-1}(j)} d_{\Psi}(\lambda_i, \lambda_j) < d_T(u, v);$$
 (13)⁶⁰⁰

$$x_{iu} + x_{iv} + x_{ju} + x_{jv} \le y_{ij} + 1$$
 $i, j \in \mathcal{L}, uv \in E(T);$ (14)₆₀₂

$$\sum_{i,j\in\mathcal{L}} y_{ij} = n - 1. \tag{15}$$

In this formulation, constraints (11) and (12) ensure that x_{605} encodes a bijection, while constraints (13) guarantee that the bijection adheres to the conditions of Theorem 1. The auxiliary variables y_{ij} in constraints (14) indicate whether a pair of leaf labels map to adjacent vertices in T, with constraint (15) ensuring the inferred migration network forms a tree. The objective function (10) facilitates the search for a solution by leveraging the relationship between population diversity and popula- 611

tion age [109, 110, 111], that suggests that more diverse populations, which are likely older, are also more probable origins of migration [38, 31, 66]. Consequently, it is more likely that such populations correspond to high-degree vertices in the tree T. Here a coefficient $\delta(i)$ represents the genetic diversity of the ith subpopulation, measured as allelic entropy averaged over all allelic positions.

2.5. Migration inference under convexity constraints

In this section, a convex label-distinctive homomorhism $\Psi \to T$ will be called *feasible*. When such homomorphism exists, T can be obtained from Ψ by a series of edge contractions, making T a *minor* of Ψ . Generally, a graph G_1 is a minor of a graph G_2 if G_1 can be obtained from G_2 by edge contractions, edge removals and node removals [112]. It is known that the problem of detecting whether a given graph is a minor of another graph is \mathcal{NP} -hard, even when input graphs are trees [113, 114]. However, minors associated with our problem satisfy a more stringent set of conditions than general graph minors: in our case, only edge contractions are allowed and, in addition, contractions of edges between labeled nodes with different labels are forbidden. We suggest that in practical settings feasible homomorphisms can be efficiently found and enumerated using dynamic programming.

The following simple property of convex homomorphisms will be useful for our subsequent analysis:

Lemma 3. Suppose that $f: \Psi \to T$ is a convex homomorphism. Then T' with is a subtree of T if and only if $\Psi' = f^{-1}(T')$ is a subtree of Ψ .

Proof. Homomorphic image of a connected subgraph is connected, as implied by the definition of a homomorphism. The converse is also true, if |T'|=1. Suppose that $|T'|\geq 2$ and the subgraph $f^{-1}(T')$ is not connected. Consider two of its connected components, Ψ'_1 and Ψ'_2 , such that the unique path P

between the nodes $\alpha \in \Psi_1'$ and $\beta \in \Psi_2'$ is shortest among all₆₄₇ such paths between different components. Then $|P| \geq 3$ and₆₄₈ $T'' = f(P \setminus \{\alpha, \beta\}) \cap T' = \emptyset$. Moreover, the vertices $f(\alpha)$ and₆₄₉ $f(\beta)$ are adjacent to vertices of T''. This leads to two distinct paths between $f(\alpha)$ and $f(\beta)$ – one in T' and another passing through T''. This contradicts the fact that T is a tree.

Next, we describe the proposed algorithmic approach. Ini-653
tially, we simplify the original phylogenetic tree Ψ by collaps-654
ing paths between leaves sharing identical labels into a single
node, a step made feasible by the convexity constraint. The resulting tree, still referred to as Ψ, may become non-binary and contains uniquely labeled leaves and possibly some labeled in-657
ternal nodes.

The algorithm performs a post-order traversal of the phylogeny Ψ and, for each node $\alpha \in V(\Psi)$, calculates a set H_{α} de-627 scribing possible homomorphisms from the subtree Ψ_{α} to sub-628 trees of T. At the root node ρ , the set H_{ρ} thus describes all 629 homomorphisms from Ψ to T. Upon completing the traversal, 630 the algorithm either concludes that no feasible homomorphism⁶⁶² 631 exists (when $H_{\rho} = \emptyset$), or initiates a pre-order traversal of Ψ .663 632 During this second traversal, it reconstructs feasible homomor-633 phisms using the information from the sets H_{α} . 665

Formally, let Λ_{α} be the set of labeled nodes in the subtree 636 Ψ_{α} . A subtree T[v, X] of T is termed an induced v-subtree if 637 it includes the vertex v, a subset X of v's neighbors, and all 638 vertices that are connected to v via paths that intersect with X 639 (Fig. 3).

For a vertex $\alpha \in V(\Psi)$, the set H_{α} consists of triples (v, X, C)called *partial homomorphism tokens* or simply *tokens*. In each

token, (i) $v \in V(T)$, (ii) $X \subseteq N_T(v)$, (iii) C is a subset of vertices

of an induced v-subtree T[v, X] such that there exists a feasible

surjective homomorphism $f: \Psi_{\alpha} \to T[v, X]$ with $f(\alpha) = v$ and $f(\Lambda_{\alpha}) = C$.

The algorithm is initialized by setting $H_{\lambda} = \{(v, \emptyset, \{v\}) :$

between the nodes $\alpha \in \Psi'_1$ and $\beta \in \Psi'_2$ is shortest among all₆₄₇ $v \in V(T)$ for all leafs λ . For an internal node α , the set H_α is such paths between different components. Then $|P| \geq 3$ and₆₄₈ constructed based on the sets from its children nodes β_1, \ldots, β_k . $T'' = f(P \setminus \{\alpha, \beta\}) \cap T' = \emptyset$. Moreover, the vertices $f(\alpha)$ and₆₄₉ The construction utilizes the following lemma:

Lemma 4. Let T[v,X] be an induced v-subtree. Then there exist a feasible surjective homomorphism $f: \Psi_{\alpha} \to T[v,X]$ with $f(\alpha) = v$ and $f(\Lambda_{\alpha}) = C$ if and only if there exist tokens $(v_1, X_1, C_1) \in H_{\beta_1}, \dots, (v_k, X_k, C_k) \in H_{\beta_k}$ satisfying the following conditions:

(a1)
$$v \sim v_i$$
 or $v = v_i$ for all $i \in \{1, \dots, k\}$;

(b1)
$$(V(T[v_i, X_i]) \setminus \{v\}) \cap (V(T[v_j, X_j]) \setminus \{v\}) = \emptyset$$
 for all $i, j \in \{1, \dots, k\}, i \neq j$;

(c1)
$$v \in V(T[v_i, X_i])$$
 if and only if $v_i = v$.

(d1) $v \notin C_i$, if α is labeled.

(e1)
$$X_i = N_T(v_i) \setminus \{v\}, if v_i \neq v.$$

(*f1*)
$$X = \{v_1, \ldots, v_k\} \setminus \{v\};$$

(g1) $C = C_1 \cup \cdots \cup C_k$, if α is unlabeled, and $C = C_1 \cup \cdots \cup C_k \cup \{v\}$, if α is labeled.

Proof. Let us prove the necessity of conditions (a1) - (g1). Suppose that there exists a feasible surjective homomorphism $f: \Psi_{\alpha} \to T[v, X]$ such that $f(\alpha) = v$. Define T_i as the image of the subtree Ψ_{β_i} under f, denoted by $f(\Psi_{\beta_i})$, and let $v_i = f(\beta_i)$. Also, let $X_i = N_T(v_i) \cap V(T_i)$ and $C_i = f(\Lambda_{\beta_i})$.

We aim to demonstrate that $T_i = T[v_i, X_i]$. According to Lemma 3, T_i is connected, which suggests that T_i must be a subgraph of $T[v_i, X_i]$. Furthermore, Lemma 3 also indicates that $f^{-1}(T[v_i, X_i])$ is connected. Given the surjectivity of f, it follows that $f^{-1}(T[v_i, X_i]) \subseteq \Psi_{\beta_i}$, leading to the conclusion that $T[v_i, X_i] \subseteq T_i$.

Consequently, the restriction of f to Ψ_{β_i} is a surjective homomorphism from Ψ_{β_i} to $T[v_i, X_i]$. Thus, the token (v_i, X_i, C_i)

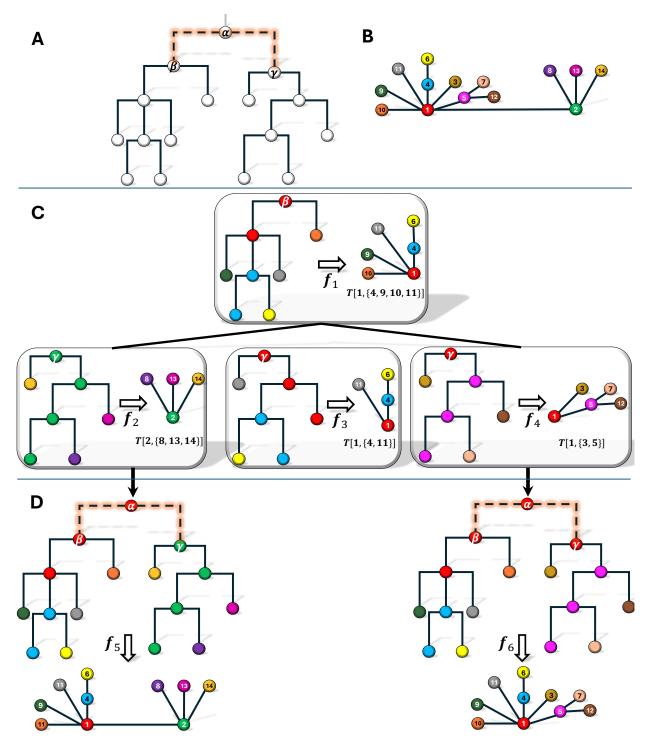


Figure 3. Overview of the Dynamic Programming Algorithm for Detecting Convex Label-Distinctive Homomorphisms. A. A phylogenetic subtree, Ψ_{α} , rooted at node α with two children, β and γ . B. Input candidate migration tree T. The goal is to produce convex homomorphisms from Ψ_{α} to induced subtrees of T from such homomorphisms for Ψ_{β} and Ψ_{γ} . C. Convex homomorphisms from Ψ_{β} (top row) and Ψ_{γ} (bottom row) to induced subtrees of T. For instance, the top figure depicts a homomorphism to an induced 1-tree $T[1, \{4, 9, 10, 11\}]$ that consists of the vertex 1, its neighbors 4, 9, 10, 11 and all vertices connected to 1 via paths that intersect these neighbors. Nodes of subtrees are colored by their homomorphic images. Homomorphisms are organized into a bipartite graph \mathcal{G} , where edges connect homomorphism pairs $f_1 f_2$ and $f_1 f_4$ that are compatible, and there is no edge between homomorphisms f_1 and f_3 , that are not compatible. D. Homomorphisms f_5 and f_6 obtained by combining compatible homomorphisms f_1, f_2 and f_1, f_4 .

- belongs to H_{β_i} . The condition (a1) follows from the homomor-679 (b1)-(g1).
- phism definition, while the convexity of f yields the conditions Conversely, suppose for each $i \in 1, ..., k$, we have feasi-

tions (a1) - (e1). We can construct a homomorphism $f: \Psi_{\alpha} \rightarrow_{711}$ T[v, X, C] by defining f as follows:

$$f(\gamma) = \begin{cases} f_i(\gamma), & \text{if } \gamma \in V(\Psi_{\beta_i}) \\ v, & \text{if } \gamma = \alpha \end{cases}$$
 (16)⁷¹⁴

Condition (a1) implies that f is a homomorphism, condition₇₁₆

(d1) ensures that labeled nodes map to distinct vertices of $T_{.717}$ and (e1) guarantees that f is surjective. It remains to show that f is convex. Suppose that $f(\gamma_1) =_{719}$ 687 $f(\gamma_2) = w$. If both γ_1 and γ_2 belong to the same subtree $\Psi_{\beta_1,720}$ 688 then the entire path $P_{\Psi}(\gamma_1, \gamma_2)$ maps to w due to the convexity of₇₂₁ 689 f_i . If, on the other hand, $\gamma_1 \in \Psi_{\beta_i}$ and $\gamma_2 \in \Psi_{\beta_i}$, then condition₇₂₂ 690 (b1) implies that w = v. Furthermore, by the condition (c1) we₇₂₃ have $v_i = v_j = v$. Therefore convexity of f_i and f_j , as well as₇₂₄ 692 the fact that $f(\alpha) = v$, imply that $f(P_{\Psi}(\gamma_1, \gamma_2)) = v$. Together,₇₂₅ these two facts prove that f is convex.

To convert Lemma 4 into an algorithm constructing the set $^{^{727}}$ of tokens for a node α using the tokens of its children, we⁷²⁸ first need to identify tokens of children that satisfy conditions 729 (a1)-(g1). This can be achieved through the following steps. 730 For each vertex $v \in V(T)$, we construct a multipartite graph⁷³¹ $\mathcal{G} = \mathcal{G}(\alpha, \nu)$ (i.e. a graph partitioned into k independent sets or 700 parts), as described below: 701

(i) The parts A_1, \ldots, A_k correspond to children of α .

702

703

704

705

706

708

- (ii) The vertices of the set A_i are tokens from the set Ψ_{β_i} satisfying the conditions (a1),(c1), (d1) and (e1).
- (iii) Two tokens from the sets A_i and A_j are adjacent whenever⁷³⁷ they satisfy the condition (b1).

In the constructed graph, sets of partial homomorphism tokens that satisfy Lemma 4 can be identified as k-vertex cliques. We employ the Bron-Kerbosch algorithm [115] to generate the-

ble homomorphisms $f_i: \Psi_{\beta_i} \to T[v_i, X_i, C_i]$ that meet condi-710 se cliques. For each identified clique, we use conditions (f1) and (g1) to construct a new token (v, X, C) for the node α .

> Additionally, for each token $(v, X, C) \in H_{\alpha}$, we maintain pointers p(v, X, C) that link to the children tokens used in its construction. These pointers are used in the subsequent phase of the algorithm, which aims to reconstruct full feasible homomorphisms f.

> During this phase, the algorithm executes a pre-order traversal of Ψ. As it progresses, it recursively assigns a specific token to each node. When a token t = (v, X, C) is assigned to node α , the algorithm sets $f(\alpha) = v$. It then retrieves tokens for t via the pointers p(t) and assigns them to children, β_1, \dots, β_k . This ensures that by the end of the traversal, each node in Ψ has been assigned a homomorphic image, completing the construction of the homomorphism f.

> In general, the set H_{ρ} at the root node ρ may include multiple tokens, each representing a feasible homomorphism f: $\Psi \to T$. When multiple feasible homomorphisms are available, the algorithm selects the one that minimizes violations of the compactness constraint (to be discussed in the next subsection). In case of ties, the selection criterion shifts to minimizing the quantity

$$D(f) = \sum_{\alpha \in \Lambda(T)} \delta(\alpha) \cdot d(f(\alpha)), \tag{17}$$

where $\delta(\alpha)$ is the number of labeled children of a node α . This approach prioritizes homomorphisms that map high-degree nodes to high-degree vertices.

The outlined method is formalized in Algorithm 2. Its efficiency can be improved by contracting sibling leaves in both Ψ and T into a single node (vertex). The algorithm maintains a count of copies of contracted nodes used by tokens, and a new token is generated from children tokens only if the total count of each leaf in the children tokens does not exceed its overall count. This modification markedly enhances the dynamic programming algorithm's runtime. The adjustments to Lemma

Algorithm 2 Convex homomorphism minimizing compactness violations

```
phylogenetic tree \Psi with the root \rho and a candidate migration tree T
                 a feasible homomorphism f:\Psi\to T or the answer that it does not exist.
1: modify \Psi by contracting paths between leafs with the same label.
2: perform a post-order traversal of \Psi.
    for every node \alpha of the post-order do
         construct a set of partial homomorphism tokens H_{\alpha}.
4:
5:
         if \alpha is a leaf then
6:
             H_{\alpha} \leftarrow \{(v, \emptyset, \{v\}) : v \in V(T)\};
7:
         end if
8:
         if \alpha is an internal node with children \beta_1,\ldots,\beta_k then
             Let H(\beta_j) = \{(v_i^J, X_i^J, C_i^J) : i = 1, ..., l_j\}, j = 1, ..., k.
9:
10:
             for v \in V(T) do
                  construct the multipartite graph \mathcal{G}(\alpha, v) as described in (i)-(iii);
11:
                 Generate the set K of k-vertex cliques of \mathcal{G}(lpha, 
u) using Bron-Kerbosch algorithm.
12:
                 for each clique \{(u_{i_1}^1, X_{i_1}^1, C_{i_1}^1), \dots, (u_{i_k}^k, X_{i_k}^k, C_{i_k}^k)\} \in K do
13:
                      Construct the sets X and C using formulas (f1) and (g1).
14:
15:
                      Set H_{\alpha} \leftarrow H_{\alpha} \cup \{(v, X, C)\} and p(v, X, C) \leftarrow p(v, X, C) \cup \{(i_1, \dots, i_k)\}.
16:
                  end for
             end for
17:
         end if
18:
19:
    end for
20:
    if H_{\rho} \neq \emptyset then
         perform a pre-order traversal of \Psi;
21:
22:
         for each token t_1, \ldots t_R \in H_{\rho} do
23:
             assign the token t_r to \rho : AS_{\rho} \leftarrow t_r.
             for every node \boldsymbol{\alpha} of the pre-order do
24:
                  f_r(\alpha) \leftarrow v, where (v, X, C) = AS_{\alpha}.
25:
                  if \alpha is an internal node with children \beta_1, \ldots, \beta_k then
26:
                     Let H(\beta_j) = \{(v_i^j, X_i^j, C_i^j) : i = 1, ..., l_j\}, j = 1, ..., k \text{ and } p(v, X, C) = (i_1, ..., i_k).
27:
                      AS_{\beta_j} \leftarrow (v_{i_i}^j, X_{i_i}^j, C_{i_i}^j), \quad j = 1, \dots, k
28:
                  end if
29:
30:
             end for
31:
         end for
         among generated homomorphisms f_1,\ldots,f_R, output the homomorphism f_r with the minimal number of compactness
    violations and, in case of ties, with the minimal D(f_r).
33: else
34:
         f does not exist
35: end if
```

4 and Algorithm 3 in Supplementary Material are straightfor-754 ward, but involve numerous minor technical details; hence a755 formal description is omitted. One particular detail, however,756 should be mentioned: if Ψ' and T' represent the leaf-contracted 757 versions of Ψ and T, respectively, and $f': \Psi' \to T'$ is a fea-758 747 sible homomorphism, there may be cases where $f(\alpha) = l$ for 748 an internal node α of Ψ and a leaf l in T' that results from the 749 contraction of leaves l_1 and l_2 . In such instances, f' can be ex-750 tended to a homomorphism $f: \Psi \to T$ by designating $f(\alpha)$ as⁷⁶¹ 751 either l_1 or l_2 . To resolve this ambiguity, the leaf corresponding⁷⁶² to the population with higher diversity is chosen. 764

Finally, it should be noted that, strictly speaking, Algorithm 2 is not polynomial, since the number of tokens for a node of Ψ theoretically can be exponential. In practical settings, however, the algorithm is extremely fast, and require split seconds to finish.

2.6. Migration inference with convexity and compactness constraints

A similar approach to the one outlined in Subsection 2.5 can be employed to identify homomorphisms that are both convex and compact. However, the dynamic programming algorithm can be further optimized by taking advantage of the specific nature

of these constraints.

As with the earlier approach, homomorphisms that are both₇₉₇ convex and compact will be referred to as feasible. Following₇₉₈ 767 a similar methodology to that used in Algorithm 2, we begin₇₉₉ 768 by contracting the paths in the phylogeny Ψ. We then construct800 769 partial homomorphism tokens similar to those used previously,801 770 with the exception that subsets C of images of labeled nodes 771 are not required. Thus, the tokens are simplified to pairs (v, X), 772 where $v \in V(T)$ and $X \subseteq N_T(v)$. A pair (v, X) is included in H_α 773 if there exists a feasible surjective homomorphism $f:\Psi_{\alpha}\to^{\mbox{\tiny 802}}$ 774 T[v,X] such that $f(\alpha)=v$, and it also satisfies the following ⁸⁰³ condition:

$$|\Lambda_{\alpha}| = |T[v, X]|. \tag{18}$$

This condition is necessary for a partial homomorphism to be 807 extendable to a full compact homomomorphism $\Psi \to T$.

The algorithm is initialized by setting $H_{\lambda} = \{(v, \emptyset) : v \in_{809} V(T)\}$ for leafs λ . For an internal node $\alpha \in V(\Psi)$, its token set 810 H_{α} is constructed from the tokens of its children $\beta_1, \dots \beta_k$ using 811 Lemma 5.

Lemma 5. $(v, X) \in H_{\alpha}$ if and only if one of the following con-813 ditions hold:

- 1) α is not labeled and there exist $w \in X$ such that $(v, X \setminus_{816} \{w\}) \in H_{\beta_1}$ and $(w, N_T(w) \setminus \{v\}) \in H_{\beta_2}$.
- 2) α is labeled, |X| = k, and there exist a permutation $(v_1, \dots, \S^1 \mathbb{P}_k)$ of elements of X such that $v \sim v_1, \dots, v_k$, and $(v_1, N_T(v_1) \setminus \S^1 \mathbb{P}_k)$ $\{v\} \in H_{\beta_1}, \dots, (v_k, N_T(v_k) \setminus \{v\}) \in H_{\beta_k}.$

Proof. We present the proof for the case where α is unlabeled; the argument for labeled α follows a similar rationale. In this case, α was not involved in path contraction, and thus k=2.

Suppose that $(v,X) \in H_{\alpha}$, i.e. $|\Lambda_{\alpha}| = |T[v,X]|$ and there exists a feasible surjective homomorphism $f: \Psi_{\alpha} \to T[v,X]$ such that $f(\alpha) = v$. Consequently, $f(\Lambda_{\alpha}) = V(T[v,X])$, and822

therefore there must be a labeled node $\gamma \in V(\Psi_{\beta_1})$ such that $f(\gamma) = v$. Given the connectivity constraint, this implies that $f(\beta_1) = v$. Meanwhile, the compactness constraint necessitates $f(\beta_2) = w \neq v$.

Following Lemma 3 and considering the connectivity constraint, it can be shown that:

$$f(\Psi_{\beta_1}) = T[v, X \setminus \{v\}] \text{ and } f(\Psi_{\beta_2}) = T[w, N(w) \setminus \{v\}]$$
 (19)

Let us now establish the second equality; the method for proving the first is analogous. Let $T_2 = f(\Psi_{\beta_2})$. We know that $v \notin V(T_2)$ and, according to 3, T_2 is connected. These observations imply that $T_2 \subseteq T[w, N(w) \setminus \{v\}]$. Conversely, Lemma 3 suggests that $f^{-1}(T[w, N(w) \setminus \{v\}])$ is connected. Given the surjectivity of f, this implies $f^{-1}(T[w, N(w) \setminus \{v\}]) \subseteq \Psi_{\beta_2}$, leading to $T[w, N(w) \setminus \{v\}] \subseteq T_2$. Thus, both $T_2 \subseteq T[w, N(w) \setminus \{v\}]$ and $T[w, N(w) \setminus \{v\}] \subseteq T_2$ are true, confirming the second equality.

So, the restrictions $f|_{\Psi_{\beta_1}}$ and $f|_{\Psi_{\beta_2}}$ are both feasible surjective homomorphisms. Additionally, $\Lambda_{\alpha} = \Lambda_{\beta_1} \cup \Lambda_{\beta_2}$, $f(\Lambda_{\beta_1}) = f(\Psi_{\beta_1}) = T[v, X \setminus \{w\}]$ and $f(\Lambda_{\beta_2}) = f(\Psi_{\beta_2}) = T[w, N_T(w) \setminus \{v\}]$, thus confirming that the equality (18) holds for the tokens $(v, X \setminus \{w\})$ and $(w, N_T(w) \setminus \{v\})$. This proves the necessity of condition 1).

To demonstrate the sufficiency of condition 1), we assume that there exist feasible surjective homomorphisms $f_1: \Psi_{\beta_1} \to T[\nu, X \setminus \{w\}]$ and $f_2: \Psi_{\beta_2} \to T[w, N_T(w) \setminus \{\nu\}]$. By defining f as follows, we can establish a combined feasible homomorphism:

$$f(\chi) = \begin{cases} f_1(\chi) & \text{if } \chi \in V(\Psi_{\beta_1}) \\ f_2(\chi) & \text{if } \chi \in V(\Psi_{\beta_2}) \\ v & \text{if } \chi = \alpha \end{cases}$$
 (20)

Lemma 5 can be directly applied to construct tokens for an unlabeled node α using the tokens from its children. For

a labeled node α , however, the process of finding a permuta-854 tion (v_1, \ldots, v_k) of the tokens is more complex. The method to achieve this is detailed in the following approach.

Lemma 5 can be straightforwardly use to construct tokens₈₅₇ of α from tokens of its children, if α is unlabeled. When α ₈₅₈ is labeled, then finding a permutation (v_1, \ldots, v_k) required by₈₅₉ Lemma 5 is more complicated and can be achieved using the₈₆₀ approach described next.

Suppose that $S = (S_1, ..., S_k)$ is a collection of sets. A vec-₈₆₂ tor $(x_1, ..., x_k)$ is termed a *transversal* of S [116] if $x_i \in S_i$ and₈₆₃ all x_i are distinct. Let now $S_i = \{u : (u, Y)\} \in H_{\beta_i}$ and $v \sim_{864} u$ }. Then the vector $(v_1, ..., v_k)$ satisfies the condition 2) of₈₆₅ Lemma 5 if and only if it is a transversal of S.

Given this, the set H_{α} can be obtained by generating all₈₆₇ transversals of S. This process involves the following steps:

836

837

838

842

843

844

845

847

851

852

- Construct a bipartite graph $\mathcal{B}(S)$ with parts I and \mathcal{J} , where I = 1, ..., k represents the set indices, and $\mathcal{J} = \bigcup_{i=1}^{k} S_i$, represents all elements in the sets. In this graph, a vertex $i \in I$ is adjacent to a vertex $j \in \mathcal{J}$ if j belongs to S_i .
- In \mathcal{B} , each transversal corresponds to a maximal matching of size k. To generate these matchings, construct a line graph $L(\mathcal{B})$, where each vertex represents an edge₈₇₆ of \mathcal{B} , and two vertices are adjacent if their corresponding edges in \mathcal{B} share a common vertex. Maximal matchings of \mathcal{B} correspond to maximal independent sets in $L(\mathcal{B})$, 878 that can be produced using the Bron–Kerbosch algorithm 879 [115].

The entire method is detailed in Supplementary Material, 882 Algorithm 3. Like Algorithm 2, efficiency can be significantly improved by contracting sibling leaves in both Ψ and T.

2.7. SMiTH: Sampling Migration Trees via Homomorphisms

The algorithms discussed can be effectively integrated into an Unlabeled Migration Sampling framework (Problem 5). It allows for the identification of homeomorphic images of a given phylogeny Ψ within a collection of candidate migration trees sampled from a given migration pattern represented by a specified tree distribution.

The obtained sample of migration trees can be directly analyzed to estimate the probabilities of specific migration routes or to obtain summary statistics and confidence intervals for derivative evolutionary parameters. It can be also synthesized into a single weighted *consensus graph*, where each edge is weighted by the number of candidate trees that support it. When the homomorphism reconstruction includes an objective function, the consensus graph is constructed from a subsample comprising the top $\kappa\%$ of trees ranked by their objective values. For applications requiring a specific output tree – such as for benchmarking and comparison with other methods described in Subsection 3.1 – the tree is determined by calculating the maximum-weight spanning tree of the consensus graph. The entire algorithmic pipeline, named SMiTH (Sampling Migration Trees via Homomorphisms), is illustrated in Figure 1.

3. Results

3.1. Simulated data

To generate synthetic data, we used FAVITES [117], a tool capable of simulating genomes, phylogenies, and migration networks under various evolutionary scenarios. Although originally designed to simulate viral outbreaks, FAVITES supports general phylogenetic and population genetics models, making it suitable for simulating migrations of heterogeneous populations besides viruses. It also should be noted that, to the best of our knowledge, specialized simulation tools for metastatic

spread with capabilities comparable to FAVITES are currently₉₁₉ not available.

We simulated the migration of a heterogeneous population₉₂₁ over a network of sites formed according to the Barabasi-Albert model [70]. This assumption can be valid for both viral [98, 118] and cancer [119, 95] spread. Migrations occur at a con-₉₂₃ stant rate along each network edge (in viral context, this corresponds to the network-based Susceptible-Infected (SI) transmission model). Within each site, phylogenies evolved under the exponential coalescent, a model previously used to simu-⁹²⁶ late intra-host [66] and intra-tumor [120] evolution. Genotypes were assumed to evolve under the GTR+Γ substitution model, and were sampled simultaneously at the end of the simulation.

In total, 275 simulated datasets were generated, encompassing 5–30 demes with 100 sequences sampled per deme.

890

891

892

893

895

900

901

902

903

904

906

907

913

914

915

916

917

918

In the first series of experiments, we sampled candidate migration trees from 3 distinct prior distributions of tree topologies and evaluated their compatibility with simulated phylogenies under three different types of constraints. The prior distributions included:

- T1) Degenerate distribution consisting of the true topology. ⁹³⁷
 Even though this scenario is unrealistic for actual migra- ⁹³⁸
 tion inference, it serves as a test to determine if migration ⁹³⁹
 links can be accurately reconstructed when the topology ⁹⁴⁰
 is known but sites need to be correctly mapped to migra- ⁹⁴¹
 tion network vertices. To avoid bias linked to correlations ⁹⁴²
 between vertex IDs and migration times, that can be po- ⁹⁴³
 tentially introduced by the simulation methods, we pro- ⁹⁴⁴
 duced multiple samples with randomly permuted vertex ⁹⁴⁵
 IDs.
- T2) Random scale-free trees produced by the preferential attachment procedure.
- T3) Uniformly distributed trees of a given size. Sampling was950

performed by generating random Prufer codes [121], integer sequences of length n-2 that uniquely define n-vertex trees.

The following types of constraints were used:

- H1) unconstrained homomorphism;
- H2) convex homomorphism minimizing the number of compactness constraint violations;
- H3) convex and compact homomorphism.

For each simulated dataset, we produced 9 samples of candidate migration trees corresponding to all combinations of conditions T1)–T3) and H1)–H3). The sample sizes ranged from 1,000 for the degenerate distribution without constraints, up to 1,000,000 for the uniform distribution with convexity and compactness constraints. This variation in sample size was necessary because stricter constraints require larger samples to ensure that a sufficient number of feasible trees are produced. We assessed the compatibility of these sampled trees with the given phylogenies under the respective constraints, using methods detailed in Subsections 2.4–2.6. Using these assessments, subsamples of compatible tree were extracted.

Sampled trees were compared with true migration trees produced by FAVITES. Individual trees were compared by measuring recall, defined as the fraction of inferred transmission edges among true transmission edges; precision, the fraction of true transmission edges among inferred transmission edges; and the *f*-score, i.e., the harmonic mean of precision and recall.

Additionally, we summarized each subsample of compatible migration trees using a *consensus graph*, where each edge is weighted by the proportion of candidate trees that support that edge [26, 34, 122]. A solution can be extracted from the consensus graph by discarding edges with support below a predefined threshold. For each graph, we estimated the area under

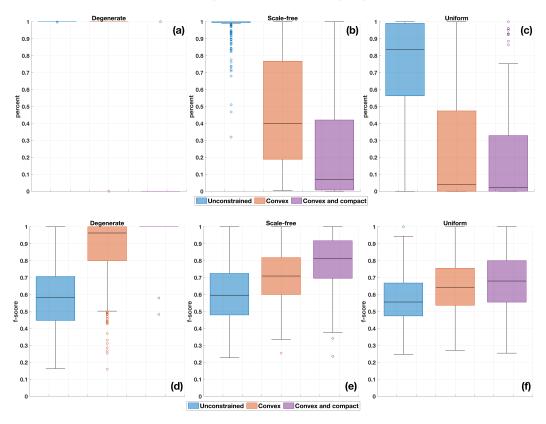


Figure 4. (a)-(c) Percent of sampled trees that are compatible with given phylogenies. (d-f) Area under the precision-recall curve (AUC) calculated by varying the support threshold.

the precision-recall curve, which was calculated by varying the support threshold. We opted for the precision-recall curve in-969 stead of the more common ROC curve due to the imbalance 970 between the classes of true and false migration edges.

We found that the relationship between phylogenetic trees₉₇₂ and migration trees is heavily influenced by structural constra-₉₇₃ ints. In the absence of constraints, there is considerable ambi-₉₇₄ guity in the possible migration histories that align with a given₉₇₅ phylogenetic tree topology, even when prior knowledge about₉₇₆ true migration tree topology is available. Notably, almost every₉₇₇ sampled scale-free network proved compatible with the given₉₇₈ phylogenies, with a median fraction of compatible trees at $\mu = 1_{979}$ (Fig. 4b). It is not entirely unexpected in light of Theorem₉₈₀ 1, which suggests that trees with a low diameter – a common₉₈₁ feature of scale-free networks – are more likely to be compati-₉₈₂ ble with a given phylogenetic tree. However, even among uni-₉₈₃ formly sampled trees, a high compatibility rate was observed₉₈₄

when no constraints are applied (median $\mu = 0.84$, Fig. 4c).

Furthermore, without constraints individual compatible trees display only marginal agreement with true migration trees. This holds not only for scale-free and uniformly sampled trees, but even for the degenerate distribution, with median f-score within the range 0.33–0.36 for all three tree priors (Supplementary Material, Fig. 8). In other words, even if the topology of a true migration tree is known, numerous labelings of that topology are compatible with the original phylogeny, most of which substantially diverge from the true labeling. Consequently, the true labeling is not immediately distinguishable among the alternatives without additional information. Combining a compatible tree subsample into a consensus graph, however, brings it closer to the true migration tree, with the median AUC values ranging from 0.56 to 0.59 for all three tree priors (see Fig. 4).

The introduction of convexity and compactness constraints significantly reduces the percentage of trees deemed compati-

ble, as can be expected (Fig. 4abc). This reduction primarily₀₁₉ eliminates incidental solutions, enhancing the alignment of the₀₂₀ trees that meet these constraints with the true migration trees₀₂₁ (Fig. 4def). In particular, for the degenerate distribution, intro₁₀₂₂ ducing constraints effectively filters out ambiguous vertex map₁₀₂₃ pings, nearly always recovering the true mapping, provided that₀₂₄ the solutions satisfying the constraints exist.

Interestingly, without constraints, the use of tree priors does₀₂₆ not enhance accuracy, as demonstrated by the lack of signif₁₀₂₇ icant differences in AUC distributions among the tree priors₀₂₈ (p = 0.083, Kruskal-Wallis test). In contrast, when constraints₀₂₉ are applied, prior knowledge of the migration tree structure be₁₀₃₀ comes beneficial, with AUCs improving as the tree prior be₁₀₃₁ comes tighter. In particular, under constraints, AUCs for the₀₃₂ scale-free prior are significantly higher than those for the uni₁₀₃₃ form prior ($p = 4.4 \cdot 10^{-6}$ and $p = 1.76 \cdot 10^{-16}$ for convex and₀₃₄ both convex and compact cases, respectively, Kruskal–Wallis₀₃₅ test). Given that true migration trees are generated by the pref₁₀₃₆ erential attachment, this suggests that under constraints, the₀₃₇ phylogeny to a certain degree reflects the properties of the un₁₀₃₈ derlying migration network.

Taken together, these observations indicate that phylogeny₀₄₀ topologies do indeed reflect underlying migration tree struc₁₀₄₁ tures, but the extent of this reflection is influenced by evolution₁₀₄₂ ary constraints. Moreover, the correspondence between phylo₁₀₄₃ genies and migration trees is primarily discernible when ana₁₀₄₄ lyzed statistically across a large sample of feasible migration₀₄₅ trees that are compatible with the phylogeny. A single compat₁₀₄₆ ible tree may be arbitrary, and thus relying on a single solution₀₄₇ may lead to misleading conclusions.

Based on those observations, we have developed a method₀₄₉ named SMiTH (Sampling MIgration Trees using Homomor₁₀₅₀ phisms) for the constrained inference of migration trees with₀₅₁ expected general properties. This method involves sampling₀₅₂

candidate migration trees from a designated random tree distribution, identifying convex homomorphisms from the given phylogeny to these sampled trees while minimizing an objective function defined by the number of compactness constraint violations, constructing a consensus graph from trees with top objective values, and ultimately inferring the final migration tree as the minimal spanning tree of this consensus graph.

We benchmarked SMiTH against several existing tools designed to infer migration networks from phylogenetic tree topologies. For the sake of fairness, tools that use dated phylogenies and/or case-specific epidemiological information were not considered. The tools selected for this comparison include Cassiopeia [21, 59], MACHINA [22], Phyloscanner [31], STraTUS [36], and TNet [44]. For MACHINA, we ran all four migration models provided by the tool and report the best result, which was achieved using the single-source seeding model. STraTUS generates a sample of migration trees rather than a single tree; thus, similarly to SMiTH, we used the minimum spanning tree of the consensus graph for benchmarking purposes.

The results of the algorithm comparison are shown in Fig. 5. It was found that both variants of SMiTH – with uniform and scale-free tree priors – allow for a statistically significant improvement over other tools ($p < 10^{-9}$, multiple comparison of f-score distributions by Kruskal–Wallis test). SMiTH is followed by Cassiopeia and STraTUS – two other sampling-based methods, whose accuracies were statistically indiscernible (p = 0.58, Kruskal–Wallis test). These tools indeed both produce samples of convex solutions, albeit using different algorithms. While STraTUS is doing it directly, while Cassiopeia's module FitchCount samples most parsimonious solutions that in our examples were almost always convex. These tools were followed by MACHINA that, similarly to our approach, imposes structural constraints on plausible migration trees by considering them as subgraphs of so-called transition patterns. However,

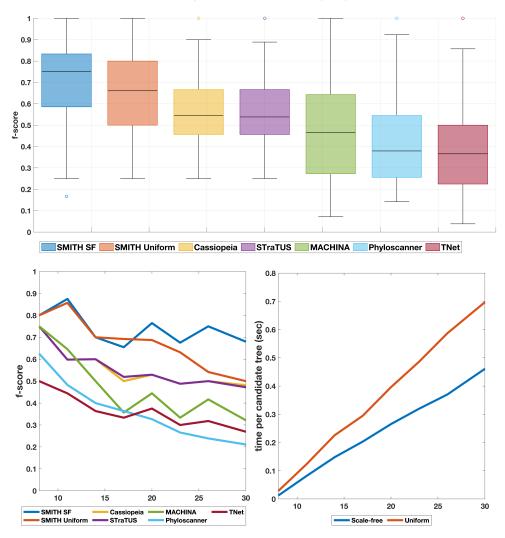


Figure 5. (a) Summary statistics of methods performance on all simulated datasets. (b) Median *f*-scores of different methods on simulated datasets with varying numbers of populations. (c) Median running times of SMiTH for construction of a constrained homomorphism for a given phylogenetic tree and candidate migration tree as a function of number of migration sites

MACHINA produces a single solution optimizing particular obtions jectives rather than summarizes a sample of such solutions; this;1066 in light of the observations described above, likely hampered its:067 performance vis-à-vis sampling-based methods. Similar rea:1068 soning can be applied to Phyloscanner, that also produces a:069 single most parsimonious solution. In addition, Phyloscanner:070 is specifically designed to make use of paraphyly, that usually:071 provides a strong signal for migration [66] when present; con:1072 sequently, its accuracy can be affected when, as in analyzed:073 test cases, the number of paraphyletic clades is limited. Fur:1074 thermore, Phyloscanner usually assumes that when the popula:1075 tions sampled from different cites are monophyletic, then they:076

all have a common source [66] - the assumption that is opposite to compactness and that seems to be not always valid. The relation between algorithms' performances is mostly stable with regard to the number of migration sites (Fig. 5b).

The median running time of our method for construction of a constrained homomorphism for a given phylogenetic tree and candidate migration tree was within 0.7 seconds for all tests and tree priors (Fig. 5c), even though theoretically our algorithms can be exponential in the worst case. It allows us, using straightforward parallelization, to produce and process samples consisting in hundreds of thousands of candidate trees in reasonable time.

3.2. Experimental viral data

1079

1080

1081

1082

1083

1084

1085

1086

1089

1090

1091

1092

1093

1094

1095

1098

1099

1100

1101

1102

1103

1104

1105

1107

1108

1109

1110

Analysis of simulated data highlights the role of structural con¹¹¹² straints in migration tree inference, particularly when paraphyly¹¹³ is limited. Interestingly, these constraints prove just as essential¹¹⁴ in scenarios with a high degree of paraphyly, albeit for different¹¹⁵ reasons. This is evidenced by the analysis of data of Hepatitis¹¹⁶ C (HCV) outbreaks, which have been considered in previous¹¹⁷ studies [38, 39, 44, 27]. The data comprises intra-host HCV¹¹⁸ populations from several outbreaks investigated by the Centers¹¹⁹ for Disease Control and Prevention, each population consisting¹²⁰ of sequences covering Hypervariable Region 1 (HVR1) of the¹²¹ HCV genome. In each outbreak, a single primary host infected¹²² all other hosts, rendering the migration tree in graph-theoretical¹²³ terms a *star*.

We analyzed two largest outbreaks involving 15 and 19 in¹¹²⁵ fected hosts. Phylogenetic trees for each outbreak were con¹¹²⁶ structed using RAxML [123]. All transmissions occurred with¹¹²⁷ in a short time frame and, as a result, intra-host populations are ¹²⁸ highly intermixed (Fig. 6a). This makes paraphylytic signal ¹²⁹ strong, but oversaturated, thus impeding its use to reconstruct true transmission history.

1131 This effect can be demonstrated by examining internal node labels generated by Fitch algorithm, that serves as a basis for several methods considered in the previous section. In the trees analyzed, 32-34% of internal nodes were assigned a single provisional label during the post-order traversal step of the dy-1136 namic programming algorithm, indicating that these labels ap-1137 pear in all most parsimonious solutions (or solutions with the minimal migration number in terms of [22]). Many of these nodes are adjacent, suggesting that the transmission links they represent will be identified by any parsimony-based sampling approach similar to those employed by existing tools [44, 12, 21]. For one outbreak, these links form a connected graph, whereas in the other, only one host does not integrate into this 1144 single connected component (see Fig. 6b). These resulting graphs are relatively dense and include not only true edges but also a significant number of false positives (see Fig. 6b). Consequently, even if all true positive edges are correctly identified using nodes with multiple Fitch labels, the f-scores would not exceed 0.54 and 0.51, respectively. Enhancing the accuracy of this approach requires the filtering out of false positive edges, achievable only through the integration of additional prior information or constraints.

In contrast, sampling unconstrained candidate transmission trees from the scale-free tree distribution produce the results that are significantly closer to true transmission histories (Fig. 6c). For consensus networks derived from these samples, areas under precision-recall curve are estimated at 0.76 and 0.70. Furthermore, f-scores of solutions obtained as minimal spanning trees of consensus networks produced from top 1% of sampled trees according to the objective (10) are 0.86 and 0.94. Comparable results -f = 0.79 and f = 0.89 – are obtained if we use top 1% of sampled trees based on the parsimony score.

3.3. Experimental cancer data

We employed SMiTH to analyze the migration history of metastatic ovarian cancer using the data published in [124]. The dataset comprised whole-genome and targeted sequencing data from samples collected at various anatomical sites, including the left ovary (LOv), the right ovary (ROv), and several metastases. In the original study, migration networks were inferred using hierarchical clustering trees and a Dollo parsimony model. Subsequent re-analysis using MACHINA [22] revealed several additional, more parsimonious migration histories.

We focused on the data from Patients 1, 3, and 7, that included the highest number of anatomical sites (7-8 sites) and that were thoroughly analyzed in [22]. We used clone trees shared by the authors of [22]. For each patient, we sampled candidate migration trees from a uniform distribution using an

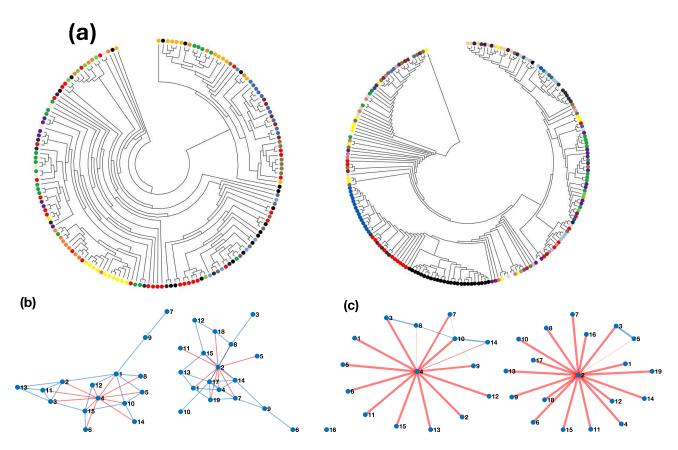


Figure 6. (a) Phylogenetic trees of HCV variants from two outbreaks. Variants sampled from different hosts are highlighted in different colors. (b) Graphs formed by edges corresponding to adjacent tree nodes with unique Fitch labels. True edges are highlighted in red. (c) Consensus networks of top 1% of sampled trees with respect to the objective (10). Edge thicknesses are proportional to their frequencies, true edges are highlighted in red.

unconstrained model. Following the methodology described in₁₅₉ [22], we resolved polytomies in clone trees to match the solu₁₁₆₀ tions reported there. In instances where the resolution of poly₁₁₆₁ tomies was ambiguous, we applied a random resolution, gen₁₁₆₂ erating a new random resolution for each sampled candidate₁₆₃ migration tree.

For Patient 1, [124] identified a complex migration history;165 designating ROv as the primary tumor site. MACHINA was166 able to find several more parsimonious histories that suggested167 either LOv or ROv as the primary tumor location, but was not168 able to distinguish between them. These histories shared parsi+169 mony scores (referred to as "migration numbers" [22]) and/or170 the same number of migration events (named "co-migration171 numbers" [22]), which left the primary tumor's location am+172

biguous. In contrast, SMiTH enabled the comparison of migration number distributions for different potential primary tumor sources (Fig. 7) rather than making the decision based on single most parsimonious solutions. Migration numbers associated with LOv and ROv were significantly lower than those of other potential sources ($p < 10^{-65}$, Mann–Whitney U test). Among these two, the lowest numbers were observed for ROv ($p < 10^{-38}$, Mann–Whitney U test), indicating a stronger statistical support for ROv as the primary tumor source.

A similar situation was observed for Patient 7. Here, [124] suggested the right uterosacral ligament (RUt) as the primary tumor location, while MACHINA identified several alternative migration histories with either the left ovary (LOv) or the right ovary (ROv) as the source, each sharing identical migration and

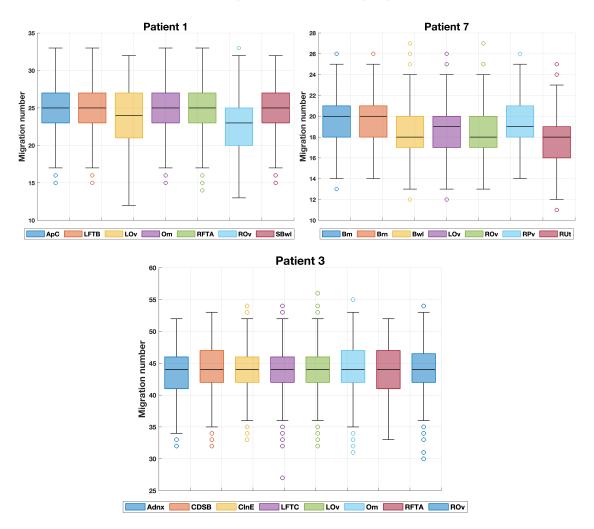


Figure 7. Distributions of migration numbers for trees with different primary tumor sites.

co-migration numbers. Based on these results, [22] argued that 186 available data provides no evidence for the assertion that the 187 primary tumor is located in the RUt as opposed to the ovaries 1188 However, SMiTH provided such statistical evidence (Fig. 7) 1189 showing that migration trees with RUt as the source generally 190 exhibited lower migration numbers ($p < 10^{-154}$, Mann–Whitney 191 U test).

Patient 3 presents a different scenario. Here, MACHINA₁₉₃ identified several migration histories with LOv or ROv as pri₁₁₉₄ mary tumor sources. There is also an alternative history with₁₉₅ the omentum (Om) as the source, which has a lower migra₁₁₉₆ tion number. The latter hypothesis appears preferable if judged₁₉₇ solely by the single most parsimonious solution. Yet, this con₁₁₉₈

clusion is not reliable due to the highly symmetric distribution of clones from different sites in the clone tree. Many clones form polytomies, offering insufficient data to clearly differentiate between the corresponding sites (e.g., all clones from LOv and LFTC are siblings, rendering these sites indistinguishable; see Supplementary Material, Fig. 9). The apparently lower migration number for Om, compared to other sites, is simply due to its representation by five clones, versus four for several other sites – a difference that could stem from sampling bias given the small number of clones involved.

SMiTH was able to capture and quantify this uncertainty. Specifically, it did not find significant differences in the distribution of migration numbers among potential primary sources, 1240

1241

with p-values ranging from 0.09 to 0.96 in pairwise Mann-1232 Whitney U tests and a p=0.65 in a joint Kruskal-Wallis233 test (Fig. 7). This provides a statistical support for suggestion234 that the existing data is insufficient to draw reliable conclusions235 about the migration history.

In total, these examples illustrate how SMiTH can be used₂₃₇ to provide statistical support for hypotheses regarding metat₁₂₃₈ static spread pathways.

4. Discussion

1200

1201

1202

1203

1204

1205

1206

1208

1209

1210

1211

1212

1213

1214

1216

1217

1218

1219

1220

1221

1222

1223

1225

1226

1227

1228

1229

1230

1231

This study is dedicated to in-depth mathematical exploration of the relationships between phylogenies and migration trees of heterogeneous genomic populations. Although it is established that phylogenetic trees impose some restrictions on migration pathways [36], the exact nature and extent of these constraints are still not well understood, despite a considerable amount of research dedicated to this problem [73, 72, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 36]. Our approach adds both depth and rigor to this area by utilizing the powerful theoretical and algorithmic framework of theory of graph homomorphisms. This framework allowed us to derive necessary and sufficient conditions for the compatibility of phylogenetic and migration trees, and to develop efficient algorithms for analyzing this compatibility through numerical experiments.

Based on our findings, we propose a general and flexible computational framework that can be used to infer migration networks under various assumptions, quantitatively assess competing hypotheses about migration dynamics, investigate the influence of phylogenies on the migration tree space, and to determine whether a potential migration history is definitively contradicted by a phylogeny or set of phylogenies.

Methodologically, our approach balances the advantages of probabilistic and parsimony methods. It incorporates scalability and the use of advanced combinatorial optimization techniques 1265

from the latter, along with the biological plausibility of the former that comes from employing appropriate prior random tree distributions.

This study aligns well with the context of previous research. Several earlier studies have conceptualized migration inference as a coloring problem [35, 34, 36], employing this framework both to develop efficient inference algorithms and to explore the structure of migration tree space. The methodology introduced in this paper advances these prior approaches by incorporating a more comprehensive mathematical model and applying more sophisticated mathematical techniques. Similarly, the concept of constraining the tree space to random trees from a specified distribution was first introduced in our earlier studies on viral transmissions [38, 27], while a related concept of restricting the tree space to subgraphs of specific migration patterns has been utilized in computational cancer genomics [22, 61]. This paper not only refines and extends these methodologies but also integrates them into a cohesive modeling and computational framework.

The proposed approach certainly has both advantages and disadvantages. We recognize that uniform sampling from the space of candidate migration trees may not be the most optimal tool for migration dynamics inference, and employing more precise optimization or sampling techniques to navigate the tree space could substantially enhance the method's accuracy and efficiency. This work lays a foundation for further theoretical and algorithmic development in this direction, equipping researchers with the tools needed to expand the use of graph homomorphism methodologies. On the other hand, uniform sampling may be more suitable for hypothesis testing and comparison. It also should be noted that a key aspect of our methodology is that it involves sampling from the space of subtrees of a transition pattern, rather than from the space of phylogeny node labelings as was done in other studies [36, 44, 12].

The former space is considerably smaller in practical scenar₁₂₉₆ ios which often makes the uniform sampling computationally, feasible.

1267

Furthermore, it is likely that migration models tailored ${\rm to}^{^{1299}}$ 1269 the specific characteristics of underlying populations could po-1270 tentially yield more accurate insights into the biological pro₇₃₀₂ 1271 cesses involved. In particular, the biological mechanisms driv¹³⁰³ 1272 ing viral and cancer migrations are certainly very different. Nonethe-1273 less, employing a unified phylogenetic approach to study highly, 306 1274 mutable populations offers several advantages. First, it decou₇₃₀₇ 1275 ples the initial phylogenetic reconstruction from its biological³⁰⁸ 1276 interpretation, thereby minimizing the risk of overfitting and 309 1277 ensuring that the results are less biased by underlying models [35]. Additionally, such methods are significantly more com₇₃₁₂ putationally efficient and scalable compared to parameter-rich³¹³ 1280 models [35, 27]. Finally, more general phylogenetic models 1281 offer greater flexibility and versatility, and usually can be read-1282 ily extended to more specific settings through the integration317 1283 of suitable priors [35]. These features make them an excellent an excellent 1284 foundation for more detailed analyses. 1285 1320

5. Acknowledgements

1286

1293

K.K. has been supported by a Georgia State University Molec₁₃₂₄ ular Basis of Disease fellowship. P.S. has been supported by³²⁵ the NSF grants 2047828 and 2212508. The authors thank the organizers of the Computational Genomics Summer Institute (UCLA, July 12 – August 4, 2023 and July 10 - August 2₁₃₂₉ 2024), where several ideas of this paper were conceived.

6. Code availability

The code developed and used in this study is available at https: 1334
//github.com/compbel/SMiTH.

References

- [1] E. Domingo, J. Sheldon, C. Perales, Viral quasispecies evolution, Microbiology and Molecular Biology Reviews 76 (2) (2012) 159–216.
- [2] R. A. Burrell, N. McGranahan, J. Bartek, C. Swanton, The causes and consequences of genetic heterogeneity in cancer evolution, Nature 501 (7467) (2013) 338–345.
- [3] E. C. Smith, The not-so-infinite malleability of rna viruses: Viral and cellular determinants of rna virus mutation rates, PLoS pathogens 13 (4) (2017) e1006254.
- [4] S. Duffy, Why are rna virus mutation rates so damn high?, PLoS biology 16 (8) (2018) e3000003.
- [5] R. Sanjuán, P. Domingo-Calap, Mechanisms of viral mutation, Cellular and molecular life sciences 73 (2016) 4433–4448.
- [6] I. P. Tomlinson, M. Novelli, W. Bodmer, The mutation rate and cancer, Proceedings of the National Academy of Sciences 93 (25) (1996) 14800–14803.
- [7] M. Greaves, Nothing in cancer makes sense except..., BMC biology 16 (1) (2018) 1–8.
- [8] W. M. Grady, Genomic instability and colon cancer, Cancer and metastasis reviews 23 (1-2) (2004) 11–27.
- [9] G. S. Charames, B. Bapat, Genomic instability and cancer, Current molecular medicine 3 (7) (2003) 589–596.
- [10] F. Utro, C. Levovitz, K. Rhrissorrakrai, L. Parida, A common methodological phylogenomics framework for intra-patient heteroplasmies to infer sars-cov-2 sublineages and tumor clones, BMC genomics 22 (5) (2021) 1–13.
- [11] T. Stadler, O. G. Pybus, M. P. Stumpf, Phylodynamics for cell biologists, Science 371 (6526) (2021) eaah6266.
- [12] P. Sashittal, M. El-Kebir, Sampling and summarizing transmission trees with multi-strain infections, Bioinformatics 36 (Supplement_1) (2020) i362-i370.
- [13] J. Z. Sanborn, J. Chung, E. Purdom, N. J. Wang, H. Kakavand, J. S. Wilmott, T. Butler, J. F. Thompson, G. J. Mann, L. E. Haydu, et al., Phylogenetic analyses of melanoma reveal complex patterns of metastatic dissemination, Proceedings of the National Academy of Sciences 112 (35) (2015) 10995–11000.
- [14] D. X. Nguyen, J. Massagué, Genetic determinants of cancer metastasis, Nature Reviews Genetics 8 (5) (2007) 341–352.
- [15] N. D. Grubaugh, J. T. Ladner, P. Lemey, O. G. Pybus, A. Rambaut, E. C. Holmes, K. G. Andersen, Tracking virus outbreaks in the twenty-first century, Nature microbiology 4 (1) (2019) 10–19.
- [16] G. L. Armstrong, D. R. MacCannell, J. Taylor, H. A. Carleton, E. B. Neuhaus, R. S. Bradbury, J. E. Posey, M. Gwinn, Pathogen genomics in

1337

1331

1332

1333

1321

- public health, New England Journal of Medicine 381 (26) (2019) 2569+382 2580. 1383
- [17] A. Black, D. R. MacCannell, T. R. Sibley, T. Bedford, Ten recommental
 dations for supporting open pathogen genomic analysis in public health
 Nature Medicine (2020) 1–10.
- [18] D. Lähnemann, J. Köster, E. Szczurek, D. J. McCarthy, S. C. Hicks₁₃₈₇
 M. D. Robinson, C. A. Vallejos, K. R. Campbell, N. Beerenwinkel₁₃₈₈
 A. Mahfouz, et al., Eleven grand challenges in single-cell data science₁₃₈₉
 Genome biology 21 (1) (2020) 1–35.
- [19] N. Navin, J. Kendall, J. Troge, P. Andrews, L. Rodgers, J. McIndoo₁₃₉₁
 K. Cook, A. Stepansky, D. Levy, D. Esposito, et al., Tumour evolution₃₉₂
 inferred by single-cell sequencing, Nature 472 (7341) (2011) 90–94. 1393
- [20] S. Knyazev, L. Hughes, P. Skums, A. Zelikovsky, Epidemiological data394
 analysis of viral quasispecies in the next-generation sequencing erat395
 Briefings in bioinformatics 22 (1) (2021) 96–108.
- J. J. Quinn, M. G. Jones, R. A. Okimoto, S. Nanjo, M. M. Chantasan
 N. Yosef, T. G. Bivona, J. S. Weissman, Single-cell lineages reveal theases
 rates, routes, and drivers of metastasis in cancer xenografts, Sciences
 371 (6532) (2021) eabc1944.
- [22] M. El-Kebir, G. Satas, B. J. Raphael, Inferring parsimonious migration401
 histories for metastatic cancers, Nature Genetics 2 (2018) 5.
- [23] Jombart, A. Cori, X. Didelot, S. Cauchemez, C. Fraser, N. Ferguson₁₄₀₃
 Bayesian reconstruction of disease outbreaks by combining epidemio₁₄₀₄
 logic and genomic data, PLoS Comput Biol 10 (1) (2014) e1003457. 1405
- [24] F. Campbell, X. Didelot, R. Fitzjohn, N. Ferguson, A. Cori, T. Jom¹⁴⁰⁶
 bart, outbreaker2: a modular platform for outbreak reconstruction, BMG⁴⁰⁷
 bioinformatics 19 (11) (2018) 1–8.
- T. Jombart, R. Eggo, P. Dodd, F. Balloux, Reconstructing disease out₁₄₀₉
 breaks from genetic data: a graph approach, Heredity 106 (2) (2011)₄₁₀
 383–390.
- 1369 [26] N. De Maio, C.-H. Wu, D. J. Wilson, Scotti: efficient reconstruction412
 1370 of transmission within outbreaks with the structured coalescent, PLoS413
 1371 computational biology 12 (9) (2016) e1005130.
- 1372 [27] P. Skums, F. Mohebbi, V. Tsyvina, P. Icer, S. Ramachandran;415

 Y. Khudyakov, Sophie: viral outbreak investigation and transmission416

 history reconstruction in a joint phylogenetic and network theory frame4417

 work, in: International Conference on Research in Computational418

 Molecular Biology, Springer, 2022, pp. 369–370.
- [28] D. Klinkenberg, J. A. Backer, X. Didelot, C. Colijn, J. Wallinga, Simul₄₂₀
 taneous inference of phylogenetic and transmission trees in infectious₄₂₁
 disease outbreaks, PLoS computational biology 13 (5) (2017) e1005495₁₄₂₂
- [29] C. J. Worby, P. D. O'Neill, T. Kypraios, J. V. Robotham, D. De Anger
 lis, E. J. Cartwright, S. J. Peacock, B. S. Cooper, Reconstructing trans

- mission trees for communicable diseases using densely sampled genetic data, The annals of applied statistics 10 (1) (2016) 395.
- [30] N. De Maio, C. J. Worby, D. J. Wilson, N. Stoesser, Bayesian reconstruction of transmission within outbreaks using genomic variants, PLoS computational biology 14 (4) (2018) e1006117.
- [31] C. Wymant, M. Hall, O. Ratmann, D. Bonsall, T. Golubchik, M. de Cesare, A. Gall, M. Cornelissen, C. Fraser, T. M. P. C. STOP-HCV Consortium, T. B. Collaboration, Phyloscanner: inferring transmission from within-and between-host pathogen genetic diversity, Molecular biology and evolution 35 (3) (2017) 719–733.
- [32] A. de Bernardi Schneider, C. T. Ford, R. Hostager, J. Williams, M. Cioce, Ü. V. Çatalyürek, J. O. Wertheim, D. Janies, Strainhub: A phylogenetic tool to construct pathogen transmission networks, Bioinformatics 36 (3) (2020) 945–947.
- [33] X. Didelot, C. Fraser, J. Gardy, C. Colijn, Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks, Molecular biology and evolution 34 (4) (2017) 997–1007.
- [34] X. Didelot, J. Gardy, C. Colijn, Bayesian inference of infectious disease transmission from whole-genome sequence data, Molecular biology and evolution 31 (7) (2014) 1869–1879.
- [35] J. Carson, M. Keeling, D. Wyllie, P. Ribeca, X. Didelot, Inference of infectious disease transmission through a relaxed bottleneck using multiple genomes per host, Molecular Biology and Evolution (2024) msad288.
- [36] M. D. Hall, C. Colijn, Transmission trees on a known pathogen phylogeny: Enumeration and sampling, Molecular biology and evolution 36 (6) (2019) 1333–1343.
- [37] S. Sledzieski, C. Zhang, I. Mandoiu, M. S. Bansal, Treefix-tp: Phylogenetic error-correction for infectious disease transmission network inference, bioRxiv (2019) 813931.
- [38] P. Skums, A. Zelikovsky, R. Singh, W. Gussler, Z. Dimitrova, S. Knyazev, I. Mandric, S. Ramachandran, D. Campo, D. Jha, et al., Quentin: reconstruction of disease transmissions from viral quasispecies genomic data, Bioinformatics 34 (1) (2017) 163–170.
- [39] O. Glebova, S. Knyazev, A. Melnyk, A. Artyomenko, Y. Khudyakov, A. Zelikovsky, P. Skums, Inference of genetic relatedness between viral quasispecies from sequencing data, BMC genomics 18 (10) (2017) 918.
- [40] S. L. Kosakovsky Pond, S. Weaver, A. J. Leigh Brown, J. O. Wertheim, Hiv-trace (transmission cluster engine): a tool for large scale molecular epidemiology of hiv-1 and other rapidly evolving pathogens, Molecular biology and evolution 35 (7) (2018) 1812–1819.
- [41] A. G. Longmire, S. Sims, I. Rytsareva, D. S. Campo, P. Skums, Z. Dimitrova, S. Ramachandran, M. Medrzycki, H. Thai, L. Ganova-Raeva,

- et al., Ghost: global hepatitis outbreak and surveillance technology,468
 BMC genomics 18 (10) (2017) 916.
- [42] E. M. Campbell, A. Boyles, A. Shankar, J. Kim, S. Knyazev, R. Cint470
 tron, W. M. Switzer, Microbetrace: retooling molecular epidemiology471
 for rapid public health response, PLoS computational biology 17 (9)472
 (2021) e1009300.

1425

1426

1431

1432

1457

1458

1459

- [43] P. Sashittal, M. El-Kebir, Sharptni: counting and sampling parsimonious₄₇₄ transmission networks under a weak bottleneck, bioRxiv (2019) 842237₁₄₇₅
- [44] S. Dhar, C. Zhang, I. Mandoiu, M. S. Bansal, Tnet: Transmission net+476
 work inference using within-host strain diversity and its application to 477
 geographical tracking of covid-19 spread, IEEE/ACM Transactions on 478
 Computational Biology and Bioinformatics.
- [45] Z. Ke, H. Vikalo, Graph-based reconstruction and analysis of disease480
 transmission networks using viral genomic data, Journal of Computa4481
 tional Biology.
- [46] R. J. Ypma, W. M. van Ballegooijen, J. Wallinga, Relating phyloge1483
 netic trees to transmission trees of infectious disease outbreaks, Genetics484
 1442
 195 (3) (2013) 1055–1062.
- [47] N. Mollentze, L. H. Nel, S. Townsend, K. Le Roux, K. Hampson, D. Ti486
 Haydon, S. Soubeyrand, A bayesian approach for inferring the dynam487
 ics of partially observed endemic infectious diseases from space-time488
 genetic data, Proceedings of the Royal Society of London B: Biological489
 Sciences 281 (1782) (2014) 20133251.
- [48] M. J. Morelli, G. Thébaud, J. Chadœuf, D. P. King, D. T. Haydon₁491
 S. Soubeyrand, A bayesian inference framework to reconstruct transmis+492
 sion trees using epidemiological and genetic data, PLoS Comput Biol493
 8 (11) (2012) e1002768.
- [49] E. M. Cottam, G. Thébaud, J. Wadsworth, J. Gloster, L. Mansley, D. J. 495
 Paton, D. P. King, D. T. Haydon, Integrating genetic and epidemiologi 496
 cal data to determine transmission pathways of foot-and-mouth disease 497
 virus, Proceedings of the Royal Society of London B: Biological Sci 498
 ences 275 (1637) (2008) 887–895.
 - [50] F. Campbell, A. Cori, N. Ferguson, T. Jombart, Bayesian inference of 500 transmission chains using timing of symptoms, pathogen genomes and 501 contact data, PLoS computational biology 15 (3) (2019) e1006930. 1502
- [51] M. Hall, M. Woolhouse, A. Rambaut, Epidemic reconstruction in a phy+503
 logenetics framework: transmission trees as partitions of the node set₁504
 PLoS computational biology 11 (12) (2015) e1004613.
- 1463 [52] M. Senghore, H. Read, P. Oza, S. Johnson, H. Passarelli-Araujo, B. Ph. B. Taylor, S. Ashley, A. Grey, A. Callendrello, R. Lee, et al., Inferring bact-507 terial transmission dynamics using deep sequencing genomic surveil+508 lance data, Nature Communications 14 (1) (2023) 6397.
 - [53] D. S. Campo, G.-L. Xia, Z. Dimitrova, Y. Lin, J. C. Forbi, L. Ganova⁴⁵¹⁰

- Raeva, L. Punkova, S. Ramachandran, H. Thai, P. Skums, et al., Accurate genetic detection of hepatitis c virus transmissions in outbreak settings, The Journal of infectious diseases 213 (6) (2016) 957–965.
- [54] J. O. Wertheim, S. L. Kosakovsky Pond, L. A. Forgione, S. R. Mehta, B. Murrell, S. Shah, D. M. Smith, K. Scheffler, L. V. Torian, Social and genetic networks of hiv-1 transmission in new york city, PLoS pathogens 13 (1) (2017) e1006000.
- [55] O. Ratmann, M. K. Grabowski, M. Hall, T. Golubchik, C. Wymant, L. Abeler-Dörner, D. Bonsall, A. Hoppe, A. L. Brown, T. de Oliveira, et al., Inferring hiv-1 transmission networks and sources of epidemic spread in africa with deep-sequence phylogenetic analysis, Nature communications 10 (1) (2019) 1–13.
- [56] Y. Zhang, C. Wymant, O. Laeyendecker, M. K. Grabowski, M. Hall, S. Hudelson, E. Piwowar-Manning, M. McCauley, T. Gamble, M. C. Hosseinipour, et al., Evaluation of phylogenetic methods for inferring the direction of human immunodeficiency virus (hiv) transmission: Hiv prevention trials network (hptn) 052, Clinical Infectious Diseases.
- [57] S. Ramachandran, H. Thai, J. C. Forbi, R. R. Galang, Z. Dimitrova, G.l. Xia, Y. Lin, L. T. Punkova, P. R. Pontones, J. Gentry, et al., A large hcv transmission network enabled a fast-growing hiv outbreak in rural indiana, 2015, EBioMedicine 37 (2018) 374–381.
- [58] E. M. Campbell, H. Jia, A. Shankar, D. Hanson, W. Luo, S. Masciotra, S. M. Owen, A. M. Oster, R. R. Galang, M. W. Spiller, et al., Detailed transmission network analysis of a large opiate-driven outbreak of hiv infection in the united states, The Journal of infectious diseases 216 (9) (2017) 1053–1062.
- [59] M. G. Jones, A. Khodaverdian, J. J. Quinn, M. M. Chan, J. A. Hussmann, R. Wang, C. Xu, J. S. Weissman, N. Yosef, Inference of single-cell phylogenies from lineage tracing data using cassiopeia, Genome biology 21 (1) (2020) 1–27.
- [60] S. Kumar, A. Chroni, K. Tamura, M. Sanderford, O. Oladeinde, V. Aly, T. Vu, S. Miura, Pathfinder: Bayesian inference of clone migration histories in cancer, Bioinformatics 36 (Supplement_2) (2020) i675–i683.
- [61] M. S. Roddur, S. Snir, M. El-Kebir, Inferring temporally consistent migration histories, in: 23rd International Workshop on Algorithms in Bioinformatics (WABI 2023), Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2023, pp. 9.1–9.22.
- [62] S. Hessey, P. Fessas, S. Zaccaria, M. Jamal-Hanjani, C. Swanton, Insights into the metastatic cascade through research autopsies, Trends in Cancer.
- [63] M. Al Bakir, A. Huebner, C. Martínez-Ruiz, K. Grigoriadis, T. B. Watkins, O. Pich, D. A. Moore, S. Veeriah, S. Ward, J. Laycock, et al., The evolution of non-small cell lung cancer metastases in tracerx, Na-

- ture (2023) 1–10.
- [64] H.-N. Chen, Y. Shu, F. Liao, X. Liao, H. Zhang, Y. Qin, Z. Wangussi
 M. Luo, Q. Liu, Z. Xue, et al., Genomic evolution and diverse models of systemic metastases in colorectal cancer, Gut 71 (2) (2022) 322–332.
- 1515 [65] R. Schwartz, A. A. Schäffer, The evolution of tumour phylogenetics1558 1516 principles and practice, Nature Reviews Genetics 18 (4) (2017) 213-4559 1517 229.
- [66] E. O. Romero-Severson, I. Bulla, T. Leitner, Phylogenetically resolvings61
 epidemiologic linkage, Proceedings of the National Academy of Sci+562
 ences (2016) 201522930.
- [67] T. Leitner, W. Fitch, The phylogenetics of known transmission histories₁₅₆₄
 The evolution of HIV. Johns Hopkins University Press, Baltimore, Md565
 (1999) 315–345.
- [68] G. E. Fischer, M. K. Schaefer, B. J. Labus, L. Sands, P. Rowley, I. At567
 Azzam, P. Armour, Y. E. Khudyakov, Y. Lin, G. Xia, et al., Hepatitis c568
 virus infections from unsafe injection practices at an endoscopy clinic in569
 las vegas, nevada, 2007–2008, Clinical Infectious Diseases 51 (3) (2010)570
 267–273.
- [69] P. Erdős, A. Rényi, et al., On the evolution of random graphs, Publ. math₁₅₇₂
 inst. hung. acad. sci 5 (1) (1960) 17–60.
- 1531 [70] A.-L. Barabási, R. Albert, Emergence of scaling in random networks, 574 science 286 (5439) (1999) 509–512.
- [71] D. J. Watts, S. H. Strogatz, Collective dynamics of 'small+576
 world'networks, nature 393 (6684) (1998) 440.
- 1535 [72] C. Colijn, J. Gardy, Phylogenetic tree shapes resolve disease transmis+578 1536 sion patterns, Evolution, medicine, and public health 2014 (1) (2014)579 1537 96–108.
- [73] G. E. Leventhal, R. Kouyos, T. Stadler, V. Von Wyl, S. Yerly, J. Böni₁₅₈₁
 C. Cellerai, T. Klimkait, H. F. Günthard, S. Bonhoeffer, Inferring epi+582
 demic contact structure from phylogenetic trees, PLoS computational583
 biology 8 (3) (2012) e1002413.
- 1542 [74] N. B. Carnegie, Effects of contact network structure on epidemic trans+585 1543 mission trees: implications for data required to estimate network struc+586 1544 ture, Statistics in medicine 37 (2) (2018) 236–248.
- [75] B. T. Grenfell, O. G. Pybus, J. R. Gog, J. L. Wood, J. M. Daly, J. A1588
 Mumford, E. C. Holmes, Unifying the epidemiological and evolutionary589
 dynamics of pathogens, science 303 (5656) (2004) 327–332.
- [76] E. M. Volz, K. Koelle, T. Bedford, Viral phylodynamics, PLoS comput-591
 tational biology 9 (3) (2013) e1002947.
- 1550 [77] R. Noble, D. Burri, C. Le Sueur, J. Lemant, Y. Viossat, J. N. Kathen, 593

 N. Beerenwinkel, Spatial structure governs the mode of tumour evolution, Nature ecology & evolution 6 (2) (2022) 207–217.

 1595
 - [78] M. A. Lewinsohn, T. Bedford, N. F. Müller, A. F. Feder, State-dependents96

- evolutionary models reveal modes of solid tumour growth, Nature Ecology & Evolution 7 (4) (2023) 581–596.
- [79] K. Robinson, N. Fyson, T. Cohen, C. Fraser, C. Colijn, How the dynamics and structure of sexual contact networks shape pathogen phylogenies, PLoS computational biology 9 (6) (2013) e1003105.
- [80] L. Villandre, D. A. Stephens, A. Labbe, H. F. Günthard, R. Kouyos, T. Stadler, S. H. C. Study, et al., Assessment of overlap of phylogenetic transmission clusters and communities in simple sexual contact networks: applications to hiv-1, PloS one 11 (2) (2016) e0148459.
- [81] D. Welch, Is network clustering detectable in transmission trees?, Viruses 3 (6) (2011) 659–676.
- [82] F. Giardina, E. O. Romero-Severson, J. Albert, T. Britton, T. Leitner, Inference of transmission network structure from hiv phylogenetic trees, PLoS computational biology 13 (1) (2017) e1005316.
- [83] R. M. McCloskey, R. H. Liang, A. F. Poon, Reconstructing contact network parameters from viral phylogenies, Virus evolution 2 (2) (2016) vew029.
- [84] D. Sankoff, Minimal mutation trees of sequences, SIAM Journal on Applied Mathematics 28 (1) (1975) 35–42.
- [85] W. Fitch, Towards defining the course of evolution: minimum change for a specified tree topology, Systematic Zoology 20 (1971) 406–416.
- [86] D. Gusfield, Integer linear programming in computational and systems biology: an entry-level text and course, Cambridge University Press, 2019.
- [87] W. M. Switzer, A. Shankar, H. Jia, S. Knyazev, F. Ambrosio, R. Kelly, H. Zheng, E. M. Campbell, R. Cintron, Y. Pan, et al., High hiv diversity, recombination, and superinfection revealed in a large outbreak among persons who inject drugs in kentucky and ohio, usa, Virus Evolution (2024) veae015.
- [88] A. J. Drummond, A. Rambaut, Beast: Bayesian evolutionary analysis by sampling trees, BMC evolutionary biology 7 (1) (2007) 214.
- [89] E. Kenah, T. Britton, M. E. Halloran, I. M. Longini Jr, Molecular infectious disease epidemiology: survival analysis and algorithms linking phylogenies to transmission trees, PLoS computational biology 12 (4) (2016) e1004869.
- [90] P. Hell, J. Nešetril, Graphs and homomorphisms, volume 28 of Oxford Lecture Series in Mathematics and its Applications, Oxford University Press, 2004.
- [91] L. Lovász, Graph minor theory, Bulletin of the American Mathematical Society 43 (1) (2006) 75–86.
- [92] N. Robertson, P. D. Seymour, Graph minors. xx. wagner's conjecture, Journal of Combinatorial Theory, Series B 92 (2) (2004) 325–357.
- [93] A. Brandstädt, V. B. Le, J. P. Spinrad, Graph classes: a survey, SIAM,

₅₉₇ 1999.

1608

1609

- [94] J. Scott, P. Kuhn, A. R. Anderson, Unifying metastasis—integrating641
 intravasation, circulation and end-organ colonization, Nature Reviews642
 Cancer 12 (7) (2012) 445–446.
- [95] P. K. Newton, J. Mason, K. Bethel, L. Bazhenova, J. Nieva, L. Nor₁₆₄₄
 ton, P. Kuhn, Spreaders and sponges define metastasis in lung cancer: a₆₄₅
 markov chain monte carlo mathematical model, Cancer research 73 (9)₆₄₆
 (2013) 2760–2769.
- [96] P. Gerlee, M. Johansson, Inferring rates of metastatic dissemination us₁₆₄₈
 ing stochastic network models, PLoS Computational Biology 15 (4)₆₄₉
 (2019) e1006868.
 - [97] F. Liljeros, C. R. Edling, L. A. N. Amaral, H. E. Stanley, Y. Åberg, Thess1 web of human sexual contacts, Nature 411 (6840) (2001) 907–908. 1652
- 1610 [98] J. O. Wertheim, A. J. Leigh Brown, N. L. Hepler, S. R. Mehta, D. D₁₆₅₃
 1611 Richman, D. M. Smith, S. L. Kosakovsky Pond, The global transmissione₅₄
 1612 network of hiv-1, The Journal of infectious diseases 209 (2) (2014) 304-4655
 1613 313.
- [99] G. J. Hughes, E. Fearnhill, D. Dunn, S. J. Lycett, A. Rambaut, A. J. L₁₆₅₇
 Brown, U. H. D. R. Collaboration, Molecular phylodynamics of the het₁₆₅₈
 erosexual hiv epidemic in the united kingdom, PLoS pathogens 5 (9)₆₅₉
 (2009) e1000590.
- [100] C. M. Romano, I. M. G. de Carvalho-Mello, L. F. Jamal, F. L. de Melo₁₆₆₁
 A. Iamarino, M. Motoki, J. R. R. Pinho, E. C. Holmes, P. M. de An₁₆₆₂
 drade Zanotto, V. Consortium, et al., Social networks shape the trans₁₆₆₃
 mission dynamics of hepatitis c virus, PLoS One 5 (6) (2010) e11170. 1664
 - [101] P. Sashittal, H. Schmidt, M. M. Chan, B. J. Raphael, Startle: a star homo+665 plasy approach for crispr-cas9 lineage tracing, bioRxiv (2022) 2022–12₁₆₆₆
- [102] C. A. Meacham, T. Duncan, The necessity of convex groups in biologi4667
 cal classification, Systematic Botany (1987) 78–90.
- [103] S. Moran, S. Snir, Convex recolorings of strings and trees: Definitions₁₆₆₉
 hardness results and algorithms, Journal of Computer and System Sci₁₆₇₀
 ences 74 (5) (2008) 850–869.
- [104] B. Bollobás, Combinatorics: set systems, hypergraphs, families of vec₄₆₇₂
 tors, and combinatorial probability, Cambridge University Press, 1986.1673
- [105] M. C. Golumbic, Algorithmic graph theory and perfect graphs, Vol. 57₁₆₇₄
 Elsevier, 2004.
- [106] M. R. Garey, R. L. Graham, D. S. Johnson, D. E. Knuth, Complexity676
 results for bandwidth minimization, SIAM Journal on Applied Matherest
 matics 34 (3) (1978) 477–495.
- 1636 [107] C. Dubey, U. Feige, W. Unger, Hardness results for approximating the 679 bandwidth, Journal of Computer and System Sciences 77 (1) (2011) 62-4680 1638 90. 1681
 - [108] A. Caprara, F. Malucelli, D. Pretolani, On bandwidth-2 graphs, Discrete682

- Applied Mathematics 117 (1-3) (2002) 1-13.
- [109] V. Puller, R. Neher, J. Albert, Estimating time of hiv-1 infection from next-generation sequence diversity, PLOS Computational Biology 13 (10) (2017) e1005775.
- [110] M. L. Russell, C. S. Fish, S. Drescher, N. A. Cassidy, P. Chanana, S. Benki-Nugent, J. Slyker, D. Mbori-Ngacha, R. Bosire, B. Richardson, et al., Using viral sequence diversity to estimate time of hiv infection in infants, PLoS Pathogens 19 (12) (2023) e1011861.
- [111] P. B. I. Baykal, J. Lara, Y. Khudyakov, A. Zelikovsky, P. Skums, Quantitative differences between intra-host hcv populations from persons with recently established and persistent infections, Virus Evolution 6 (2) (2021) year 103
- [112] R. Diestel, Graph theory, volume 173 of, Graduate texts in mathematics (2012) 7.
- [113] P. Kilpeläinen, H. Mannila, Ordered and unordered tree inclusion, SIAM Journal on Computing 24 (2) (1995) 340–356.
- [114] J. Matoušek, R. Thomas, On the complexity of finding iso-and other morphisms for partial k-trees, Discrete Mathematics 108 (1-3) (1992) 343–364.
- [115] C. Bron, J. Kerbosch, Algorithm 457: finding all cliques of an undirected graph, Communications of the ACM 16 (9) (1973) 575–577.
- [116] L. Lovász, M. D. Plummer, Matching theory, Vol. 367, American Mathematical Soc., 2009.
- [117] N. Moshiri, M. Ragonnet-Cronin, J. O. Wertheim, S. Mirarab, Favites: simultaneous simulation of transmission networks, phylogenetic trees and sequences, Bioinformatics 35 (11) (2019) 1852–1861.
- [118] A. J. L. Brown, S. J. Lycett, L. Weinert, G. J. Hughes, E. Fearnhill, D. T. Dunn, Transmission network parameters estimated from hiv sequences for a nationwide epidemic, Journal of Infectious Diseases (2011) jir550.
- [119] S. P. Castillo, R. A. Rebolledo, M. Arim, M. E. Hochberg, P. A. Marquet, Metastatic cells exploit their stoichiometric niche in the network of cancer ecosystems, Science Advances 9 (50) (2023) eadi7902.
- [120] D. Posada, Cellcoal: coalescent simulation of single-cell sequencing samples, Molecular biology and evolution 37 (5) (2020) 1535–1542.
- [121] S. Caminiti, I. Finocchi, R. Petreschi, On coding labeled trees, Theoretical computer science 382 (2) (2007) 97–108.
- [122] C. Wymant, F. Blanquart, T. Golubchik, A. Gall, M. Bakker, D. Bezemer, N. J. Croucher, M. Hall, M. Hillebregt, S. H. Ong, et al., Easy and accurate reconstruction of whole hiv genomes from short-read sequence data with shiver, Virus evolution 4 (1) (2018) vey007.
- [123] A. Stamatakis, Raxml version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies, Bioinformatics 30 (9) (2014) 1312– 1313.

[124] A. McPherson, A. Roth, E. Laks, T. Masud, A. Bashashati, A. W. Zhang,
 G. Ha, J. Biele, D. Yap, A. Wan, et al., Divergent modes of clonal spread
 and intraperitoneal mixing in high-grade serous ovarian cancer, Nature
 genetics 48 (7) (2016) 758–767.

7. Supplementary Material

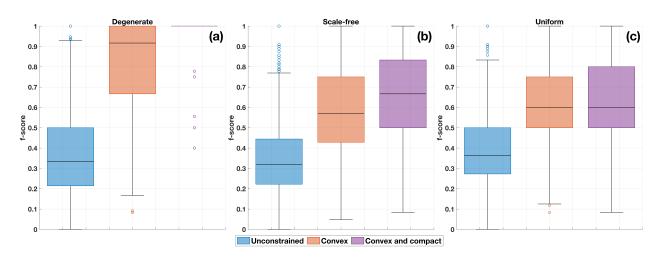


Figure 8. Comparison of sampled compatible trees with true trees under different constraints.

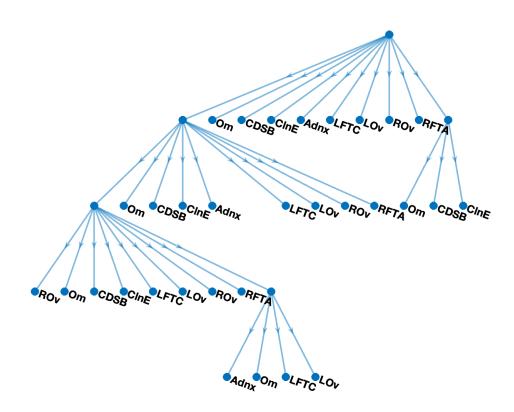


Figure 9. Clone tree for Patient 3

Algorithm 3 Homomorphism under convexity and sampling parsimony constraints

```
phylogenetic tree \Psi with the root \rho and a candidate migration tree T
     Output: a feasible homomorphism f:\Psi\to T or the answer that it does not exist.
1: modify \Psi by contracting paths between leafs with the same label.
2: perform a post-order traversal of the tree \Psi.
3: for every node \alpha of the post-order do
         construct a set of partial homomorphism tokens H_{\alpha}.
4:
         if \alpha is a leaf then
5:
             H_{\alpha} \leftarrow \{(v, \emptyset) : v \in V(T)\};
6:
7:
         end if
8:
         if \alpha is an unlabeled internal node with children \beta_1 and \beta_2 then
9:
             for all (v, X) \in H_{\beta_1} and (u, Y) \in H_{\beta_2} do
10:
                  if v \sim u, u \notin X and N_T(u) = Y \cup \{v\} then
11:
                      Add (v, X \cup \{u\}) to H_{\alpha}
12:
                  end if
                  if v \sim u, v \notin Y and N_T(v) = X \cup \{u\} then
13:
                      Add (u, Y \cup \{v\}) to H_{\alpha}
14:
                  end if
15:
              end for
16:
         end if
17:
18:
         if \alpha is a labeled internal node with children \beta_1, \ldots, \beta_k then
              S_i \leftarrow \{v : (v, X) \in H_{\beta_i} \text{ and } |X| = \deg(v) - 1\}, i = 1, k;
19:
20:
              Generate the set Tr of transversals of the set system (S_1, \ldots, S_k);
21:
              for transversals V = (v_1, \dots, v_k) \in Tr do
22:
                  if there exists v \sim v_1, \dots, v_k then
23:
                      Add (v, V) to H_{\alpha}
                  end if
24:
              end for
25:
         end if
26:
27: end for
28: if H_0 \neq \emptyset then
         perform a pre-order traversal of \Psi;
29:
         for each token t_1, \ldots t_R \in H_{\varrho} do
30:
31:
              assign the token t_r to \rho: AS_{\rho} \leftarrow t_r.
32:
              for every node \alpha of the pre-order do
33:
                  f(\alpha) \leftarrow v, where (v, X) = AS_{\alpha}.
                  if \alpha is an unlabeled internal node with children eta_1 and eta_2 then
34:
                      \text{select a vertex } u \in X \text{ such that } (v, X \setminus \{u\}) \in H_{\beta_1} \text{ and } (u, N_T(u) \setminus \{v\}) \in H_{\beta_2} \,.
35:
36:
                      AS_{\beta_1} \leftarrow (v, X \setminus \{u\}) \text{ and } AS_{\beta_2} \leftarrow (u, N_T(u) \setminus \{v\})
37:
38:
                  if \alpha is a labeled internal node with children \beta_1, \ldots, \beta_k then
39:
                      AS_{\beta_i} \leftarrow (v_i, N_T(v_i) \setminus \{v\}), i = 1, \dots, k, \text{ where } (v, \{v_1, \dots, v_k\}) = AS_{\alpha}
40:
41:
              end for
42:
              among generated homomorphisms f_1, \ldots, f_R, output the homomorphism f_r with the minimal D(f_r).
43:
44: else
45:
         f does not exist
46: end if
```