Global Alignment of Molecular Sequences via Ancestral State Reconstruction

Alexandr Andoni, Constantinos Daskalakis, Avinatan Hassidim, Sebastien Roch[§]

December 14, 2009

Abstract

Molecular phylogenetic techniques do not generally account for such common evolutionary events as site insertions and deletions (known as indels). Instead tree building algorithms and ancestral state inference procedures typically rely on substitution-only models of sequence evolution. In practice these methods are extended beyond this simplified setting with the use of heuristics that produce global alignments of the input sequences an important problem which has no rigorous model-based solution. In this paper we open a new direction on this topic by considering a version of the multiple sequence alignment in the context of stochastic indel models. More precisely, we introduce the following *trace reconstruction problem on a tree* (TRPT): a binary sequence is broadcast through a tree channel where we allow substitutions, deletions, and insertions; we seek to reconstruct the original sequence from the sequences received at the leaves of the tree. We give a recursive procedure for this problem with strong reconstruction guarantees at low mutation rates, providing also an alignment of the sequences at the leaves of the tree. The TRPT problem without indels has been studied in previous work (Mossel 2004, Daskalakis et al. 2006) as a bootstrapping step towards obtaining information-theoretically optimal phylogenetic reconstruction methods. The present work sets up a framework for extending these works to evolutionary models with indels.

In the TRPT problem we begin with a random sequence x_1, \ldots, x_k at the root of a *d*-ary tree. If vertex v has the sequence y_1, \ldots, y_{k_v} , then each one of its *d* children will have a sequence which is generated from y_1, \ldots, y_{k_v} by flipping three biased coins for each bit. The first coin has probability p_s for Heads, and determines whether this bit will be substituted or not. The second coin has probability p_d , and determines whether this bit will be deleted, and the third coin has probability p_i and determines whether a new random bit will be inserted. The input to the procedure is the sequences of the *n* leaves of the tree, as well as the tree structure (but not the sequences of the inner vertices) and the goal is to reconstruct an approximation to the sequence of the root (the DNA of the ancestral father). For every $\epsilon > 0$, we present a deterministic algorithm which outputs an approximation of x_1, \ldots, x_k if $p_i + p_d < O(1/k^{2/3} \log n)$ and $(1 - 2p_s)^2 > O(d^{-1} \log d)$.

To our knowledge, this is the first rigorous trace reconstruction result on a tree in the presence of indels.

^{*}CSAIL, MIT

[†]CSAIL, MIT. costis@mit.edu. Part of this work was done while the author was a postdoctoral researcher at Microsoft Research. [‡]MIT

[§]Department of Mathematics, UCLA. Part of this work was done while the author was a postdoctoral researcher at Microsoft Research.

1 Introduction

Trace reconstruction on a star. In the "trace reconstruction problem" (TRP) [Lev01a, Lev01b, BKKM04, KM05, HMPW08, VS08], a random binary string X of length k generates an i.i.d. collection of traces Y_1, \ldots, Y_n that are identical to X except for random *mutations* which consist in *indels*, i.e., the deletion of an old site or the insertion of a new site between existing sites, and *substitutions*, i.e., the flipping of the state at an existing site¹. (In keeping with biological terminology, we refer to the components or positions of a string as *sites*.) The goal is to reconstruct efficiently the original string with high probability from as few random traces as possible.

An important motivation for this problem is the reconstruction of ancestral DNA sequences in computation biology [BKKM04, KM05]. One can think of X as a gene in an (extinct) ancestor species 0. Through speciation, the ancestor 0 gives rise to a large number of descendants $1, \ldots, n$ and gene X evolves independently through the action of mutations into sequences Y_1, \ldots, Y_n respectively. Inferring the sequence X of an ancient gene from extant descendant copies Y_1, \ldots, Y_n is a standard problem in evolutionary biology [Tho04]. The inference of X typically requires the solution of an auxiliary problem, the *multiple sequence alignment problem* (which is an important problem in its own right in computational biology): site t_i of sequence Y_i and site t_j of sequence Y_j are said to be *homologous* (in this simplified TRP setting) if they descend from a common site t of X only through substitutions; in the multiple sequence alignment problem, we seek roughly to uncover the homology relation between Y_1, \ldots, Y_n . Once homologous sites have been identified, the original sequence X can be estimated, for instance, by site-wise majority.

The TRP as defined above is an *idealized* version of the ancestral sequence reconstruction problem in one important aspect. It ignores the actual phylogenetic relationship between species 1, ..., n. A *phylogeny* is a (typically, binary) tree relating a group of species. The leaves of the tree correspond to extant species. Internal nodes can be thought of as extinct ancestors. In particular the root of the tree represents the most recent common ancestor of all species in the tree. Following paths from the root to the leaves, each bifurcation indicates a speciation event whereby two new species are created from a parent. An excellent introduction to phylogenetics is [SS03].

A standard assumption in computational phylogenetics is that genetic information evolves from the root to the leaves according to a Markov model on the tree. Hence, the stochastic model used in trace reconstruction can be seen as a special case where the phylogeny is *star-shaped*. (The substitution model used in trace reconstruction is known in biology as the Cavender-Farris-Neyman (CFN) [Cav78, Far73, Ney71] model.) It may seem that a star is a good first approximation for the evolution of DNA sequences. However extensive work on the so-called "reconstruction problem" in theoretical computer science and statistical physics has highlighted the importance of taking into account the full tree model in analyzing the reconstruction of ancestral sequences.

The "reconstruction problem." In the "reconstruction problem" (RP), we have a single site which evolves through substitutions only from the root to the leaves of a tree. In the most basic setup which we will consider here, the tree is *d*-ary and each edge is an independent symmetric indel-free channel where the probability of a substitution is a constant $p_s > 0$. The goal is to reconstruct the state at the root given the vector of states at the leaves. More generally, one can consider a sequence of length *k* at the root where each site evolves independently according to the Markov process above. Denote by *n* the number of leaves in the tree. The RP has attracted much attention in the theoretical computer science literature due to its deep connections to computational phylogenetics [Mos03, Mos04, DMR06, Roc08] and statistical physics [Mos98, EKPS00, Mos01, MP03, MSW04, JM04, BKMP05, BCMR06, GM07, BVVD07, Sly09a, Sly09c]. See e.g. [Roc07, Sly09b] for background.

Unlike the star case, the RP on a tree exhibits an interesting thresholding effect: on the one hand, information is lost at an exponential rate along each path from the root; on the other hand, the number of paths grows exponentially with the number of levels. When the substitution probability is low, the latter "wins" and vice versa. This "phase transition" has been thoroughly analyzed in the theoretical computer science and mathematical physics literature— although much remains to be understood. More formally, we say that the RP is *solvable* when the correlation between the root and the leaves persists no matter how large the tree is. Note that unlike the TRP we do not require

¹One can also consider the case where X is arbitrary rather than random. We will not discuss this problem here.

high-probability reconstruction in this case as it is not information-theoretically achievable for d constant—simply consider the information lost on the first level below the root. Moreover the "number of traces" is irrelevant here as it is governed by the depth of the tree and the solvability notion implies nontrivial correlation for any depth. When the RP is unsolvable, the correlation decays to 0 for large trees. The results of [BRZ95, EKPS00, Iof96, BKMP05, MSW04, BCMR06] show that for the CFN model, if $p_s < p^*$, then the RP is solvable, where $d(1-2p^*)^2 = 1$. This is the so-called *Kesten-Stigum* bound [KS66]. If, on the other hand, $p_s > p^*$, then the RP is *unsolvable*. Moreover in this case, the correlation between the root state and any function of the character states at the leaves decays as $n^{-\Omega(1)}$. The positive result above is obtained by taking a majority vote over the leaf states.

Like the TRP, the RP is only an *idealized* version of the ancestral sequence reconstruction problem: it ignores the presence of indels. In other words, the RP assumes that the multiple sequence alignment problem has been solved perfectly. This is in fact a long-standing assumption in evolutionary biology where one typically preprocesses sequence data by running it through a multiple sequence alignment heuristic and then one only has to model the substitution process. This simplification has come under attack in the biology literature, where it has been argued that alignment procedures often create systematic biases that affect analysis [LG08, WSH08]. Much empirical work has been devoted to the proper joint estimation of alignments and phylogenies [TKF91, TKF92, Met03, MLH04, SR06, RE08, LG08, LRN⁺09].

Our results. We make progress in this recent new direction by analyzing the RP in the presence of indels—which we also refer to as the TRP on a tree (TRPT). We consider a *d*-ary tree where each edge is an independent channel with substitution probability p_s , deletion probability p_d , and insertion probability p_i (see Section 1.1 for a precise statement of the model). The root sequence has length *k* and is assumed to be uniform in $\{0, 1\}^k$. As in the standard RP, we drop the requirement of high-probability reconstruction and seek instead a reconstructed sequence that has correlation with the true root sequence uniformly bounded in the depth.

We give an efficient recursive procedure which solves the TRPT for $p_s > 0$ a small enough constant (strictly below, albeit close, to the Kesten-Stigum bound) and p_d , $p_i = O(k^{-2/3} \log^{-1} n)$. As a by-product of our analysis we also obtain a partial global alignment of the sequences at the leaves. Our method provides a framework for separating the indel process from the substitution process by identifying well-preserved subsequences which then serve as markers for alignment and reconstruction (see Section 1.2 for a high-level description of our techniques). As far as we are aware, our results are the first rigorous results for this problem.

Results on the RP have been used in previous work to advance the state of the art in rigorous phylogenetic tree reconstruction methods [Mos04, DMR06, MHR08, Roc08]. A central component in these methods is to solve the RP on a partially reconstructed phylogeny to obtain sequence information that is "close" to the evolutionary past; then this sequence information is used to obtain further structural information about the phylogeny. The whole phylogeny is built by alternating these steps. Our method sets up a framework for extending these techniques beyond substitution-only models. Partial results of this type will be given in the full version of the paper.

Related work. Much work has been devoted to the trace reconstruction problem on a star [Lev01a, Lev01b, BKKM04, KM05, HMPW08, VS08]. In particular, in [HMPW08], it was shown that, when there are only deletions, it is possible to tolerate a small constant deletion rate using poly(k) traces. For a different range of parameters, Viswanathan and Swaminathan [VS08] showed that, under constant substitution probability and $O(1/\log k)$ indel probability, $O(\log k)$ traces suffice. Both results assume that the root sequence X is uniformly random.

The multiple sequence alignment problem as a combinatorial optimization problem (finding the best alignment under some pairwise scoring function) is known to be NP-hard [WJ94, Eli06]. Most heuristics used in practice, such as CLUSTAL [HS88], T-Coffee [NHH00], MAFFT [KMKM02], and MUSCLE [Edg04], use the idea of a guide tree, that is, they first construct a very rough phylogenetic tree from the data (using edit distance as a measure of evolutionary distance), and then recursively construct local alignments produced by "aligning alignments." Our work can be thought of as the first attempt to analyze rigorously this type of procedure.

Finally, our work is tangentially related to the study of edit distance. Edit distance and pattern matching in random environments has been studied, e.g., by [Nav01, NBYST, AK08].

1.1 Definitions

We now define our basic model of sequence evolution.

Definition 1.1 (Model of sequence evolution) Let $T_H^{(d)}$ be the d-ary tree with H levels and $n = d^H$ leaves. For simplicity, we assume throughout that d is odd. We consider the following model of evolution on $T_H^{(d)}$. The sequence at the root of $T_H^{(d)}$ has length k and is drawn uniformly at random over $\{0,1\}^k$. Along each edge of the tree, each site (or position) undergoes the following mutations independently of the other sites:

- Substitution. The site state is flipped with probability $p_s > 0$.
- **Deletion.** The site is deleted with probability $p_d > 0$.
- Insertion. A new site is created to the right of the current site with probability $p_i > 0$. The state of this new site is uniform $\{0, 1\}$.

These operations occur independently of each other. The last two are called indels. We let $p_{id} = p_i + p_d$ and $\theta_s = 1 - 2p_s$. The parameters p_s, p_d, p_i may depend on k and n, where n is the number of leaves.

Remark 1.2 For convenience, our model of insertion is intentionally simplistic. In the biology literature, related continuous-time Markov models are instead used for this kind of process [TKF91, TKF92, Met03, MLH04, RE08, DR09]. It should be possible to extend our results to such generalizations by proper modifications to the algorithm.

1.2 Results

Statement of results. Our main result is the following. Denote by $X = x_1, \ldots, x_k$ a binary uniform sequence of length k. Run the evolutionary process on $T_H^{(d)}$ with root sequence X and let Y_1, \ldots, Y_n be the sequences obtained at the leaves, where $Y_i = y_1^i, \ldots, y_k^i$.

Theorem 1.3 (Main result) For all $\chi > 0$ and $\beta = O(d^{-1} \log d)$, there is $\Phi, \Phi', \Phi'' > 0$ such that the following holds for d large enough. There is a polynomial-time algorithm \mathbb{A} with access to Y_1, \ldots, Y_n such that for all

$$(1 - 2p_{\rm s})^2 > \frac{\Phi \log d}{d}, \qquad p_{\rm i} + p_{\rm d} < \frac{\Phi'}{k^{2/3} \log n}, \qquad \Phi'' \log^3 n < k < \text{poly}(n),$$

the algorithm A outputs a binary sequence \widehat{X} which satisfies the following with probability at least $1 - \chi$:

- 1. $\hat{X} = \hat{x}_1, \dots, \hat{x}_k$ has length k.
- 2. For all j = 1, ..., k, $\mathbb{P}[\hat{x}_j = x_j] > 1 \beta$.

Remark 1.4 Notice that we assume that the (leaf-labelled) tree and and the sequence length of the root are known. The requirement that the sequence length is known is not crucial. We adopt it for simplicity in the presentation.

Remark 1.5 In fact, we prove a stronger result which allows $\chi = o(1)$ and shows that the "agreement" between \hat{X} and X "dominates" an i.i.d. sequence. See Lemma B.1 and Section 5.2.

Proof sketch. We give a brief proof sketch. As discussed previously, in the presence of indels the reconstruction of ancestral sequences requires the solution of the *multiple sequence alignment* problem. However, in addition to being computationally intractable, global alignment through the optimization of a pairwise scoring function may create biases and correlations that are hard to quantify. Therefore, we require a more probabilistic approach. From a purely information-theoretic point of view the pairwise alignment of sequences that are far apart in the tree is

difficult. A natural solution to this problem is instead to perform *local* alignments and ancestral reconstructions, and recurse our way up the tree.

This *recursive* approach raises its own set of issues. Consider a parent node and its d children. It may be easy to perform a local alignment of the children's sequences and derive a good approximation to the parent sequence (for example, through site-wise majority). Note however that, to allow a recursion of this procedure all the way to the root, we have to provide strong guarantees about the probabilistic behavior of our local ancestral reconstruction. As is the case for global alignment, a careless alignment procedure creates biases and correlations that are hard to control. For instance, it is tempting to treat misaligned sites as independent unbiased noise but this idea presents difficulties:

Consider a site j of the parent sequence and suppose that for this site we have succeeded in aligning all but two of the children, say 1 and 2. Let $x_{j_i}^i$ denote the site in the *i*'th child which was used to estimate the *j*'th site. By the independence assumption on the root sequence and the inserted sites, $x_{j_1}^1$ and $x_{j_2}^2$ are uniform and independent of $(x_{j_i}^i)_{i=3}^d$. However, $x_{j_1}^1$ and $x_{j_2}^2$ may originate from the same neighboring site of the parent sequence and therefore are themselves correlated.

Quantifying the effect of this type of correlation appears to be nontrivial.

Instead, we use an *adversarial* approach to local ancestral reconstruction. That is, we treat the misaligned sites as being controlled by an adversary who seeks to flip the reconstructed value. This comes at a cost: it produces an asymmetry in our ancestral reconstruction. Although the RP is well-studied in the symmetric noise case, much remains to be understood in the asymmetric case. In particular, obtaining tight results in terms of substitution probability here may not be possible as the critical threshold of the RP may be hard to identify. We do however provide a tailored analysis of the particular instance of the RP by recursive majority obtained through this adversarial approach and we obtain results that are close to the known threshold for the symmetric case. Unlike the standard RP, the reconstruction error is not i.i.d. but we show instead that it "dominates" an i.i.d. noise. (See Section 4.2 for a definition.) This turns out to be enough for a well-controlled recursion. We first define a local alignment procedure which has a fair success probability (independent of n). However, applying this alignment procedure is somewhat robust in the sense that even if one of the d inputs to the reconstruction procedure is faulty, it still has a good probability of success.

As for our local alignment procedure, we adopt an *anchor* approach. Anchors were also used by [KM05, HMPW08]—although in a quite different way. We imagine a partition of every node's sequence into islands of length $O(k^{1/3})$. (The precise choice of the island length comes from a trade-off between the length and the number of islands in bounding the "bad" events below—see the proof of Lemma 3.3.) At the beginning of each island we have an anchor of length $O(\log n)$. Through this partition of the sequences in islands and anchors we aim to guarantee the following. Given a specific father node v, with fair probability 1) all the anchors in the children nodes are indel-free; and 2) for all parent islands, almost all of the corresponding children islands—those that do not satisfy these properties—are treated as controlled by an adversary. We show that Conditions 1) and 2) are sufficient to guarantee that: the anchors of all islands can be aligned with high probability and single indel events between anchors can be identified. This allows a local alignment of all islands with at most one "bad" child per island and is enough to perform a successful adversarial recursive majority vote as described above. The bound on the maximum indel probability sustained by our reconstruction algorithm comes from satisfying Conditions 1) and 2) above.

Notation. For a sequence $X = x_1, \ldots, x_k$, we let $X[i : j] = x_i, \ldots, x_j$. We use the expression "with high probability (w.h.p.)" to mean "with probability at least 1 - 1/poly(n)" where the polynomial in n can be made of arbitrarily high degree (by choosing the appropriate constants large enough). We denote by Bin(n, p) a random variable with binomial distribution of parameters n, p. For two random variables X, Y we denote by $X \sim Y$ the equality in distribution.

Organization. The rest of the paper is organized as follows. We describe the algorithm in Section 2. The proof of our main result is divided into two sections. In Section 3, we prove a series of high-probability claims about the

evolutionary process. Then, conditioning on these claims, we provide a deterministic analysis of the correctness of the algorithm in Section 5.2. All proofs are in the Appendix.

2 Description of the Algorithm

In this section we describe our algorithm for TRPT. Our algorithm is recursive, proceeding from the leaves of the tree to the root. We describe the recursive step applied to a non-leaf node of the tree.

Recursive Setup—Our Goal. For our discussion in this section, let us consider a non-leaf node v with d children, denoted u_i for $i \in [d]$. For notational convenience, we drop the index u and denote its children by $1, \ldots, d$. Our goal for the recursive step of the algorithm is to reconstruct the sequence at the node v given the sequences of the children. Denote the sites of the father by $X_0 = x_1^0, \ldots, x_{k_0}^0$, and the sites of the *i*'th child by $X_i = x_1^i, \ldots, x_{k_i}^i$. During the reconstruction process, we do not have access to the children's sequences, but rather to reconstructed sequences denoted by $\hat{X}_i = \hat{x}_1^i, \ldots, \hat{x}_k^i$.

Let us consider the following partition of the sequence of v into subsequences, called *islands*. Of course our algorithm doesn't have access to the sequence at v during the recursive step of the algorithm. We define the partition as a means to describe our algorithm: The sites of v are partitioned into *islands* of length $\ell = k^{1/3}$ (except for the last one which is possibly shorter). Denote by $N_0 = \lceil k_0/\ell \rceil$ the number of islands in v. Each island starts with an *anchor* of a bits. That is, the islands are the bitstrings $X_0[1 : \ell], X_0[\ell + 1 : 2\ell], \ldots$ and the anchors are the bitstrings $X_0[1 : a], X_0[\ell + 1 : \ell + a], \ldots$

Our algorithm tries to identify for each island $X_0[(i-1)\ell + 1 : i\ell]$ the substrings of each of the *d* children that correspond to this island (i.e., contain the sites of the island), called "child islands." We do so iteratively for $i = 1 \dots N_0$. We use the islands that did not have indels for sequence reconstruction, using the substitution-only model. Some islands will have indels however. This leads to two "modes of failure": one invalidates the entire (parent) node, and the other invalidates only an island of a child. More specifically, a node becomes invalidated (i.e., useless) when indels are not evenly distributed, that is: when an indel occured in an anchor, or two (or more) indels occured in a specific island over all *d* children. This is a rare event. Barring this event, we expect that each island suffers only at most one indel over all children. The island (of a child) that has exactly one indel is invalidated (second mode of failure), and is thus deemed useless for reconstruction purposes. As long as the parent node is not invalidated, each island will have at least d - 2 non-invalidated children islands (one additional island is potentially lost to a child node that may have been invalidated at an earlier stage).

Even when the algorithm identifies that a child island has an indel somewhere, the island is not ignored. The algorithm still needs to compute the length of the island in order to know the start of the next island in this child. For this purpose, we use the anchor of the next island and match it to the corresponding anchors of the other (non-invalidated) child islands. In fact the same procedure lets us detect which of the child islands are invalidated.

More formally, we define d functions $f_i : \{1, \ldots, k_0\} \to \{1, \ldots, k_i\} \cup \{\dagger\}$, where f_i takes the sites of v to the corresponding sites of the *i*'th child or to the special symbol \dagger if the site was deleted. Note that for each *i*, f_i is monotone, when ignoring sites which are mapped to \dagger . For $t = \ell r$, let $s_i(r) = f_i(t+1) - (t+1)$ denote the displacement of the site corresponding to the $(t+1)^{\text{st}}$ site of the parent, in the *i*th child. By convention, we take $s_i(0) = 0$. If there is no indel between $t = \ell r$ and $t' = \ell r'$ then $s_i(r) = s_i(r')$. Note that, in the specific case of one indel operation in the island, we have that $|s_i(r) - s_i(r')| = 1$.

Algorithm. Our algorithm estimates the values of $s_i(r)$ and uses these estimates to match the starting positions of the islands in the children. The full algorithm is given in Figure 1 in the Appendix. We use the following additional notation. For $x \in \{0, 1\}$ we let $\langle x \rangle = 2x - 1$. Then, for two $\{0, 1\}^m$ -sequences $Y = y_1, \ldots, y_m$ and $Z = z_1, \ldots, z_m$, we define their (empirical) correlation as

$$\operatorname{Corr}(Y, Z) = \frac{1}{m} \sum_{j=1}^{m} \langle y_j \rangle \langle z_j \rangle.$$

Note that $y \mapsto \langle y \rangle$ maps 1 to 1 and 0 to -1. One can think of Corr(Y,Z) as a form of normalized centered

Hamming distance between Y and Z. In particular, a large value of Corr(Y, Z) implies that Y and Z tend to agree. We will use the following threshold (which will be justified in Section 5.1)

$$\gamma = ((1 - \delta)(1 - 2p_{\rm s})^2 - 4\beta),$$

where δ is chosen so that

$$(1-\delta)(1-2p_{\rm s})^2 - 8\beta > \delta + 8\beta,$$

where again $\beta = O(d^{-1} \log d)$.

3 Analyzing the Indel Process

We define $a \ge C \log n$ and $\alpha \le \varepsilon/d < 1$, for constants C, ε to be determined later. We require $a < k^{1/3} < poly(n)$.² We assume that the indel probability per site satisfies

$$p_{\rm id} = \frac{\alpha}{4dk^{2/3}a} = O\left(\frac{1}{k^{2/3}\log n}\right)$$

Throughout, we denote the tree by T = (V, E).

3.1 Bound on the Sequence Length

As the indel probability is defined per site, longer sequences suffer more indel operations than shorter ones. We begin by bounding the effect of this process. We show that with high probability the lengths of all sequences are roughly equal.

Lemma 3.1 (Bound on sequence length) For all $\zeta > 0$ (small), there exists C' > 0 (large) so that for all u in V, we have

$$k_v \in [\underline{k}, \overline{k}] \equiv [(1 - \zeta)k, (1 + \zeta)k],$$

with high probability given $k \ge C' \log^3 n$. We denote this event by \mathcal{L} .

3.2 Existence of a Dense Stable Subtree

In this section, we show that with probability close to 1 there exists a dense subtree of T with a "good indel structure," as defined below. Our algorithm will try to identify this subtree and perform reconstruction on it, as described in Section 4.

Indel structure of a node. Recall that $\ell = k^{1/3}$.

Definition 3.2 (Indel structure) For a node (parent) v, we say that v is radioactive if one of the following events happen:

- 1. Event \mathcal{B}_1 : Node v has a child u such that when evolving from v to u an indel operation occurred in at least one of the sites which are located in an anchor.
- 2. Event \mathcal{B}_2 : There is an island I and two children u, u', such that an indel occurred in I in the transition from v to u and in the transition from v to u'.
- 3. Event \mathcal{B}_3 : There is an island I and a child u, such that two indel operations (or more) happened in I in the transition from v to u.

²A variant of the algorithm where the anchors have length $O(\log k)$ also works when $k \gg n$.

Otherwise the node v is stable. By definition, the leaves of T are stable. A subtree of T is stable if all of its nodes are stable.

Lemma 3.3 (Bound on radioactivity) For all $\alpha > 0$, there exists a choice of $\zeta > 0$ small enough in Lemma 3.1 such that conditioning on the event \mathcal{L} occuring: any vertex v is radioactive with probability at most α .

As a corollary we obtain the following.

Lemma 3.4 (Existence of a dense stable subtree) For all $\chi > 0$, there is of $\zeta > 0$ small enough in Lemma 3.1 such that, conditioning on the event \mathcal{L} occuring, with probability at least $1 - \chi$, the root of T is the father of a (d-1)-ary stable subtree of T. We denote this event by S.

4 A Stylized Reconstruction Process

In this subsection, we lay out the basic lemmas that we need to analyze our ancestral reconstruction method. We do this by way of describing a hypothetical sequence reconstruction process performed on the stable tree defined by the indel process (see Lemma 3.4). We analyze this reconstruction process (assuming that the radioactive nodes and the islands with indels are controlled by an adversary) and show in Lemma 4.5 that the process gives strong reconstruction guarantees. Then we argue in Section 5 that our algorithm performs at least as well as the reconstruction process against the adversary described in this section. Throughout this section we suppose that a stable tree exists and is given to us, together with the "orbit" of every site of the sequence at the root of the tree (see function F below). However, we are given no information about the substitution process.

Let $v \in V$ and assume v is the root of a (d-1)-ary stable subtree $T^* = (V^*, E^*)$ of T. (We make the stable subtree below v into a (d-1)-ary tree by potentially removing arbitrary nodes from it, at random.) Let $u \in V^*$. For each island I in u, at most one child u' of u in T^* contains an indel in which case it contains exactly one indel. We say that such an I is a corrupted island of u'. The basic intuition behind our analysis is that, provided the alignment on T^* is performed correctly (which we defer to Section 5.2), the ancestral reconstruction step of our algorithm is a recursive majority procedure against an adversary which controls the corrupted islands and the radioactive nodes (as well as all their descendants). Below we analyze this adversarial process.

Recursive majority. We begin with a formal definition of recursive majority. Let Maj : $\{0, 1, \sharp\}^d \to \{0, 1\}$ be the function that returns the majority value over non- \sharp values, and flips an unbiased coin in case of a tie (including the all- \sharp vector). Let $n_0 = d^{H_0}$ be the number of leaves in T below v. Consider the following recursive function of $z = (z_1, z_2, \ldots, z_{n_0}) \in \{0, 1, \sharp\}$: Maj⁰ $(z_1) = z_1$, and

$$Maj^{j}(z_{1},...,z_{d^{j}}) = Maj(Maj^{j-1}(z_{1},...,z_{d^{(j-1)}}),...,Maj^{j-1}(z_{d^{j}-d^{(j-1)}+1},...,z_{d^{j}})),$$

for all $j = 1, ..., H_0$. Then, $\operatorname{Maj}^{H_0}(z)$ is the *d*-wise recursive majority of *z*.

Let $X_0 = x_1^0, \ldots, x_{k_0}^0$ be the sequence at v. For $u \in V^*$ and $t = 1, \ldots, k_0$, we denote by $F_u(t)$ the position of site x_t^0 in u or \dagger if the site has been deleted on the path to u. We say that $C_{u,t}$ holds if $F_u(t)$ is in a corrupted island of u. Let Path(u, v) be the set of nodes on the path between u and v.

Definition 4.1 (Gateway node) A node u is a gateway for site t if:

- 1. $F_u(t) \neq \dagger$; and
- 2. For all $u' \in Path(u, v) \{v\}$, $C_{u',t}$ does not hold.

We let $T_t^{**} = (V_t^{**}, E_t^{**})$ be the subtree of T^* containing all gateway nodes for t. By construction, T_t^{**} is at least (d-2)-ary and for convenience we remove arbitrary nodes, at random, to make it exactly (d-2)-ary. Notice that, for $t, t' \in [1 : k_0]$, the subtrees T_t^{**} and $T_{t'}^{**}$ are random and correlated. However, they are independent of the substitution process.

We will show in Section 5.2 that the reconstructed sequence produced by our method at v "dominates" (see below) the following reconstruction process. Let $L_v = u_1, \ldots, u_{n_0}$ be the leaves below v ordered according to a planar realization of the subtree below v. Denote by $X_i = x_1^i, \ldots, x_{k_i}^i$ the sequence at u_i . For $t = 1, \ldots, k_0$, let L_t^{**} be the leaves of T_t^{**} . We define the following auxiliary sequences: for $u_i \in L_v$, we let $\Xi_i = \xi_1^i, \ldots, \xi_{k_i}^i$ where for $t = 1, \ldots, k_0$

$$\xi_t^i = \begin{cases} x_{F_{u_i}(t)}^i & \text{if } u_i \in L_t^* \\ 1 - x_t^0 & \text{o.w.} \end{cases}$$

In words, ξ_t^i is the descendant of x_t^0 if u_i is a gateway to t and is the opposite of the value x_t^0 otherwise. Because of the monotonicity of recursive majority, the latter choice is in some sense the "worst adversary" (ignoring correlations between sites—we will come back to this point later). We then define a reconstructed sequence at v as $\hat{\Xi}_0 = \hat{\xi}_1^0, \ldots, \hat{\xi}_{k_0}^0$ where for $t = 1, \ldots, k_0$

$$\hat{\xi}_t^0 = \operatorname{Maj}^{H_0}(\xi_t^1, \dots, \xi_t^{n_0}).$$

We now analyze the accuracy of this (hypothetical) estimator—which we refer to as the *adversarial reconstruction* of X_0 . We show in Section 5.2 that our actual estimator is at least as good as $\widehat{\Xi}_0$ w.h.p.

4.1 Recursive Majority Against an Adversary

To analyze the performance of the adversarial reconstruction $\widehat{\Xi}_0$, we consider the following stylized process.

Definition 4.2 (Recursive Majority Against an Adversary) We consider the following process:

- 1. Run the evolutionary process on $T_{H_0}^{(d-2)}$ at one position only starting with root state 0 without indels, that is, taking $p_{id} = 0$.
- 2. Then complete $T_{H_0}^{(d-2)}$ into $T_{H_0}^{(d)}$ and associate to each additional node the state 1.
- 3. Let $R_{H_0}^{(d)}$ be the random variable in $\{0, 1\}$ obtained by running recursive majority on the leaf states obtained above.

We call this process the *recursive majority against an adversary on* $T_{H_0}^{(d)}$. We show the following.

Lemma 4.3 (Accuracy of recursive majority) For all $\beta > 0$, there exists a constant C'' > 0 such that taking

$$\theta_{\rm s}^2 > \frac{C'' \log d}{d},$$

and d large enough, then the probability that the recursive majority against an adversary on $T_{H_0}^{(d)}$ correctly reconstructs root state 0 is at least $1 - \beta$ uniformly in H_0 . In comparison, note that the Kesten-Stigum bound for binary symmetric channels on d-ary trees is $\theta^2 > d^{-1}$ [KS66, Hig77].

As a corollary of Lemma 4.3, we have the following.

Definition 4.4 (Bernoulli sequence) For q > 0 and $m \in \mathbb{N}$, the (q,m)-Bernoulli sequence is the product distribution on $\{0,1\}^m$ such that each position is 1 independently with probability 1 - q. We denote by $B_{q,m}$ the corresponding random variable.

Lemma 4.5 (Subsequence reconstruction) Assume v is the root of a (d-1)-ary stable subtree. For all $\beta > 0$, choosing C'' > 0 as in Lemma 4.3 is such that the following holds for d large enough. For $t, m \in \{1, \ldots, k_0\}$, let $\Lambda = (\lambda_1, \ldots, \lambda_m)$ be the agreement vector between the $\widehat{\Xi}_0[t+1:t+m]$ and $X_0[t+1:t+m]$, that is, $\lambda_i = 1$ if recursive majority correctly reconstructs position i. Then there is $0 \le \beta' \le \beta$ such that $\Lambda \sim B_{\beta',m}$. (Here, β' may depend on H_0 but β does not.)

4.2 Stochastic Domination and Correlation

In our discussion so far we have assumed that a stable tree exists and is given to us, together with the function F. This allowed us to define the stylized recursive majority process against an adversary (Definition 4.2), for which we established strong reconstruction guarantees (Lemmas 4.3 and 4.5). In reality, we have no access to the stable tree. We are going to construct it recursively from the leaves towards the root. At the same time we will align sequences, discover corrupted islands, and reconstruct sequences of internal nodes. The stylized recursive majority process will be used to provide a "lower bound" on the actual reconstruction process. The notion of "lower bound" that is of interest to us is captured by *stochastic domination*, which we proceed to define formally.

Definition 4.6 (Stochastic domination) Let X, Y be two random variables in $\{0, 1\}^m$. We say that Y stochastically dominates X, denoted $X \leq Y$, if there is a joint random variable (\tilde{X}, \tilde{Y}) such that the marginals satisfy $X \sim \tilde{X}$ and $Y \sim \tilde{Y}$ and moreover $\mathbb{P}[\tilde{X} \leq \tilde{Y}] = 1$.

Note that in the definition above X and Y may (typically) live in different probability spaces. Then, the joint variable (\tilde{X}, \tilde{Y}) is a coupled version of X and Y. In our case, X is the adversarial recursive process whereas Y is the actual reconstruction performed by the algorithm. We now show how to use this property for correlation estimation.

Correlation. The analysis of the previous section guarantees that the sequences output by the adversarial reconstruction process are well correlated with the true sequences. But if we are only going to use the adversarial process as a lower bound for the true reconstruction process, it is important to establish that stochastic domination preserves correlation. In preparing the ground for such a claim let us establish an important property of the adversarial process. Let T_u and T_v be the two disjoint copies of $T_h^{(d)}$ rooted at the nodes u and v respectively, and let $X = x_1, x_2, \ldots, x_m \in \{0, 1\}^m$ and $Y = y_1, y_2, \ldots, y_m \in \{0, 1\}^m$ be sequences at the nodes u and v. Assume that u and v are the roots of (d-1)-ary stable subtrees. Let $\widehat{X}' = \widehat{x}'_1, \widehat{x}'_2, \ldots, \widehat{x}'_m \in \{0, 1\}^m$ and $\widehat{Y}' = \widehat{y}'_1, \widehat{y}'_2, \ldots, \widehat{y}'_m \in \{0, 1\}^m$ be the reconstructions of X and Y obtained by the adversarial reconstruction process. Let $\Lambda = \lambda_1, \ldots, \lambda_m$ and $\Theta = \theta_1, \ldots, \theta_m$ be the resulting agreement vectors. We show the following:

Lemma 4.7 (Concentration of bias) Let β', β be as in Lemma 4.5. Then, with probability at least $1 - e^{-\Omega(m\beta^2)}$ the following are satisfied

$$\left|\frac{1}{m}\sum_{i=1}^{m}\langle\lambda_i\rangle\langle\theta_i\rangle - (1-2\beta')^2\right| \le \frac{1}{2}\beta;$$
$$\left|\frac{1}{m}\sum_{i=1}^{m}\mathbb{1}_{\langle\lambda_i\rangle=-1} - \beta'\right| \le \frac{1}{2}\beta; \qquad \left|\frac{1}{m}\sum_{i=1}^{m}\mathbb{1}_{\langle\theta_i\rangle=-1} - \beta'\right| \le \frac{1}{2}\beta.$$

We use the previous lemma to argue that stochastic domination does not affect our correlation computations.

Lemma 4.8 (Correlation bound) Let $\hat{X}, \hat{Y} \in \{0, 1\}^m$ be random strings defined on the same probability space as \hat{X}' and \hat{Y}' . Denote by Z (resp. W) the agreement vectors of \hat{X} (resp. \hat{Y}) with X (resp. Y). Assume that $\Lambda \leq Z$ and $\Theta \leq W$ with probability 1, where Λ and Θ are the agreement vectors of \hat{X}' and \hat{Y}' with X and Y as explained above. Then,

$$|\operatorname{Corr}(X,Y) - \operatorname{Corr}(\widehat{X},\widehat{Y})| \le 1 - \frac{1}{m} \sum_{i=1}^{m} (\langle \lambda_i \rangle \langle \theta_i \rangle - \mathbb{1}_{\langle \lambda_i \rangle = -1} - \mathbb{1}_{\langle \theta_i \rangle = -1}),$$

with probability 1. Furthermore, conditioned on the conclusions of Lemma 4.7, we have, with probability 1:

$$|\operatorname{Corr}(X,Y) - \operatorname{Corr}(\widehat{X},\widehat{Y})| \le 8\beta.$$

5 Analyzing the True Reconstruction Process

We provide the proof of Theorem 1.3. In Section 5.1, we show that, if a stable subtree exists, the adversarial reconstructions of aligned anchors exhibit strong correlation signal, while misaligned anchors exhibit weak signal. This holds true for sequences that stochastically dominate the adversarial reconstructions. We use this property to complete the analysis of our reconstruction method in Section 5.2.

5.1 Anchor Alignment

Consider a parent v that is stable. Let i, j be two children with sequences $X_i = x_1^i, \ldots, x_{k_i}^i$ and $X_j = x_1^j, \ldots, x_{k_j}^j$. Let $t = \ell r$ and consider the following subsequences (of length a) at i and j

$$\mathscr{A}_{r}^{i} = x^{i}[t + s_{i}(r) + 1 : t + s_{i}(r) + a], \text{ and } \mathscr{A}_{r}^{j} = x^{j}[t + s_{j}(r) + 1 : t + s_{j}(r) + a]$$

These are related (but not identical) to the definition of anchors in the algorithm of Section 2. In particular, note that by definition \mathscr{A}_r^i and \mathscr{A}_r^j are always aligned, in the sense that they correspond to the same subsequence of v. Consider also the following subsequences

$$\mathscr{D}_r^j = x^j[t+s_j(r):t+s_j(r)+a-1]$$
 and $\mathscr{I}_r^j = x^j[t+s_j(r)+2:t+s_j(r)+a+1].$

These are the one-site shifted subsequences for j. The following lemma bounds the correlation between these strings. More precisely, we show that \mathscr{A}_r^i is always significantly more correlated to its aligned brother \mathscr{A}_r^j than to the misaligned ones \mathscr{D}_r^j and \mathscr{I}_r^j . This follows from the fact that the misaligned subsequences are sitewise independent.

Lemma 5.1 (Anchor correlations) For all $\delta > 0$ such that $(1 - \delta)(1 - 2p_s)^2 - 8\beta > \delta + 8\beta$, there is C > 0 large enough so that with $a = C \log n$, the following hold:

1. Aligned anchors.
$$\mathbb{P}\left[\operatorname{Corr}(\mathscr{A}_r^i,\mathscr{A}_r^j) > (1-\delta)(1-2p_s)^2\right] > 1 - \exp\left(-\Omega(a)\right) = 1 - 1/\operatorname{poly}(n).$$

2. *Misaligned anchors.* $\mathbb{P}\left[\operatorname{Corr}(\mathscr{A}_r^i, \mathscr{D}_r^j) < \delta\right] > 1 - \exp\left(-\Omega(a)\right) = 1 - 1/\operatorname{poly}(n)$, and similarly for \mathscr{I}_r^j .

We denote by $A_{i,j,r}$ the above events and their symmetric counterparts under $i \leftrightarrow j$.

Lemma 5.2 (Anchor correlations: Reconstructed version) Let $\hat{X}_i = (\hat{x}_i^i)_{i=1}^{k_i}$ and $\hat{X}_j = (\hat{x}_i^j)_{i=1}^{k_j}$ dominate the adversarial reconstructions \hat{X}'_i and \hat{X}'_j of X_i and X_j , as defined in Lemma 4.8. Let $\hat{\mathscr{A}}_r^i = \hat{x}^i[t + s_i(r) + 1 : t + s_i(r) + a]$ and similarly for all other possibilities $\hat{\mathscr{A}} \leftrightarrow \hat{\mathscr{D}}, \hat{\mathscr{I}}$ and/or $i \leftrightarrow j$. Denote by $\mathcal{B}_{i,j,r}$ the event that the conclusions of Lemma 4.7 hold for \hat{X}'_i and \hat{X}'_j over all pairs of intervals involving $[t + s_i(r) : t + s_i(r) + a - 1]$, $[t + s_i(r) + 1 : t + s_i(r) + a]$, and $[t + s_i(r) + 2 : t + s_i(r) + a + 1]$, with $i \leftrightarrow j$ as necessary. Then, conditioned on $\mathcal{B}_{i,j,r}$ we have

$$\operatorname{Corr}(\hat{\mathscr{A}}_{r}^{i}, \hat{\mathscr{A}}_{r}^{j}) > (1 - \delta)(1 - 2p_{s})^{2} - 8\beta,$$
$$\operatorname{Corr}(\hat{\mathscr{A}}_{r}^{i}, \hat{\mathscr{D}}_{r}^{j}) < \delta + 8\beta, \qquad \operatorname{Corr}(\hat{\mathscr{A}}_{r}^{i}, \hat{\mathscr{A}}_{r}^{j}) < \delta + 8\beta,$$

as well as their symmetric counterparts under $i \leftrightarrow j$.

5.2 **Proof of Correctness**

We show that our recursive procedure reconstructs the desired sequence at the root of the tree whenever a collection of good events occurs. Recall the definitions of the events \mathcal{L} , \mathcal{S} , $\mathcal{B}_{i,j,r}$, $\mathcal{A}_{i,j,r}$ from Lemmas 3.1, 3.4, 5.1 and 5.2.³

³Event \mathcal{L} guarantees that there is no big variance in the nodes' sequence lengths; event \mathcal{S} guarantees that a stable (d-1)-ary subtree exists; the events $\mathcal{B}_{i,j,r}$ guarantee that the adversarial reconstruction process is successful, also in preserving correlations between sequences of nodes; and the events $\mathcal{A}_{i,j,r}$ guarantee that aligned anchors (across sequences of a node's children) exhibit strong correlation signal, while misaligned anchors give weak correlation signal.

Conditioning on \mathcal{L} and \mathcal{S} , denote by $T^* = (V^*, E^*)$ the stable (d-1)-ary subtree of T. Then, for all $v \in V^*$, all pairs of children i, j of v in T^* , and all $r = 1, \ldots, \overline{k}/\ell$, we condition on the events $\mathcal{B}_{i,j,r}$ and $\mathcal{A}_{i,j,r}$. Note that having conditioned on \mathcal{L} there is only a polynomial number of such events, since all sequence lengths are bounded by \overline{k} . (If $r \cdot \ell$ is larger than a node's sequence length we assume that the corresponding events are vacuously satisfied.) Finally recall that, conditioning on \mathcal{L} , the event \mathcal{S} occurs with probability $1 - \chi$ and all other events occur with high probability. We denote the collection of events by \mathcal{E} .

Conditioning on \mathcal{E} , the proof of correctness of the algorithm follows from a bottom-up induction. The gist of the argument is the following. Suppose that at a recursive step of the algorithm we have reconstructed sequences for all children of a node v, which are strongly correlated with the true sequences (in the sense of dominating the corresponding adversarial reconstructions). Having conditioned on the events $\mathcal{A}_{i,j,r}$ and $\mathcal{B}_{i,j,r}$, it follows then that the correct alignments of anchors exhibit strong correlation signal while the incorrect alignments weak correlation signal. Hence, our correlation tests between anchors discover the corrupted islands and do the anchor alignments correctly (at least for all nodes lying inside the stable tree). Hence the shift functions \hat{s}_i 's are correctly inferred, and the reconstruction of v's sequence can be shown to dominate the corresponding adversarial reconstruction. The complete proof details are given in App endix D.

References

- [AK08] A. Andoni and R. Krauthgamer. The smoothed complexity of edit distance. Lecture Notes in Computer Science, 5125:357–369, 2008.
- [BCMR06] Christian Borgs, Jennifer T. Chayes, Elchanan Mossel, and Sébastien Roch. The Kesten-Stigum reconstruction bound is tight for roughly symmetric binary channels. In *FOCS*, pages 518–530, 2006.
- [BKKM04] Tuğkan Batu, Sampath Kannan, Sanjeev Khanna, and Andrew McGregor. Reconstructing strings from random traces. In SODA '04: Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms, pages 910–918, Philadelphia, PA, USA, 2004. Society for Industrial and Applied Mathematics.
- [BKMP05] N. Berger, C. Kenyon, E. Mossel, and Y. Peres. Glauber dynamics on trees and hyperbolic graphs. *Probab. Theory Rel.*, 131(3):311–340, 2005. Extended abstract by Kenyon, Mossel and Peres appeared in proceedings of 42nd IEEE Symposium on Foundations of Computer Science (FOCS) 2001, 568–578.
- [BRZ95] P. M. Bleher, J. Ruiz, and V. A. Zagrebnov. On the purity of the limiting Gibbs state for the Ising model on the Bethe lattice. *J. Statist. Phys.*, 79(1-2):473–482, 1995.
- [BVVD07] N. Bhatnagar, J. Vera, E. Vigoda, and Weitz D. Reconstruction for colorings on trees. Preprint available at arxiv.org/abs/0711.3664, 2007.
- [Cav78] J. A. Cavender. Taxonomy with confidence. *Math. Biosci.*, 40(3-4), 1978.
- [DMR06] Constantinos Daskalakis, Elchanan Mossel, and Sébastien Roch. Optimal phylogenetic reconstruction. In STOC'06: Proceedings of the 38th Annual ACM Symposium on Theory of Computing, pages 159–168, New York, 2006. ACM.
- [DR09] Constantinos Daskalakis and Sebastien Roch. Alignment-free phylogenetic reconstruction. Preprint, 2009.
- [Edg04] Robert C. Edgar. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucl. Acids Res.*, 32(5):1792–1797, 2004.

- [EKPS00] W. S. Evans, C. Kenyon, Y. Peres, and L. J. Schulman. Broadcasting on trees and the Ising model. Ann. Appl. Probab., 10(2):410–433, 2000.
- [Eli06] Isaac Elias. Settling the intractability of multiple alignment. *Journal of Computational Biology*, 13(7):1323–1339, 2006. PMID: 17037961.
- [Far73] J. S. Farris. A probability model for inferring evolutionary trees. *Syst. Zool.*, 22(4):250–256, 1973.
- [GM07] Antoine Gerschenfeld and Andrea Montanari. Reconstruction for models on random graphs. In FOCS '07: Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science, pages 194–204, Washington, DC, USA, 2007. IEEE Computer Society.
- [Hig77] Y. Higuchi. Remarks on the limiting Gibbs states on a (d + 1)-tree. *Publ. Res. Inst. Math. Sci.*, 13(2):335–348, 1977.
- [HMPW08] Thomas Holenstein, Michael Mitzenmacher, Rina Panigrahy, and Udi Wieder. Trace reconstruction with constant deletion probability and related results. In *SODA '08: Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 389–398, Philadelphia, PA, USA, 2008. Society for Industrial and Applied Mathematics.
- [HS88] D. G. Higgins and P. M. Sharp. Clustal: a package for performing multiple sequence alignment on a microcomputer. *Gene*, 73(1):237–244, 1988.
- [Iof96] D. Ioffe. On the extremality of the disordered state for the Ising model on the Bethe lattice. *Lett. Math. Phys.*, 37(2):137–143, 1996.
- [JM04] S. Janson and E. Mossel. Robust reconstruction on trees is determined by the second eigenvalue. *Ann. Probab.*, 32:2630–2649, 2004.
- [KM05] S. Kannan and A. McGregor. More on reconstructing strings from random traces: insertions and deletions. In *Proceedings of ISIT*, pages 297–301, 2005.
- [KMKM02] Kazutaka Katoh, Kazuharu Misawa, Kei-ichi Kuma, and Takashi Miyata. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucl. Acids Res.*, 30(14):3059–3066, 2002.
- [KS66] H. Kesten and B. P. Stigum. Additional limit theorems for indecomposable multidimensional Galton-Watson processes. Ann. Math. Statist., 37:1463–1481, 1966.
- [Lev01a] Vladimir I. Levenshtein. Efficient reconstruction of sequences. *IEEE Transactions on Information Theory*, 47(1):2–22, 2001.
- [Lev01b] Vladimir I. Levenshtein. Efficient reconstruction of sequences from their subsequences or supersequences. J. Comb. Theory Ser. A, 93(2):310–332, 2001.
- [LG08] Ari Loytynoja and Nick Goldman. Phylogeny-Aware Gap Placement Prevents Errors in Sequence Alignment and Evolutionary Analysis. *Science*, 320(5883):1632–1635, 2008.
- [LRN⁺09] Kevin Liu, Sindhu Raghavan, Serita Nelesen, C. Randal Linder, and Tandy Warnow. Rapid and Accurate Large-Scale Coestimation of Sequence Alignments and Phylogenetic Trees. *Science*, 324(5934):1561–1564, 2009.
- [Met03] Dirk Metzler. Statistical alignment based on fragment insertion and deletion models. *Bioinformatics*, 19(4):490–499, 2003.

- [MHR08] Radu Mihaescu, Cameron Hill, and Satish Rao. Fast phylogeny reconstruction through learning of ancestral sequences. *CoRR*, abs/0812.1587, 2008.
- [MLH04] I. Miklos, G. A. Lunter, and I. Holmes. A "Long Indel" Model For Evolutionary Sequence Alignment. Mol Biol Evol, 21(3):529–540, 2004.
- [Mos98] E. Mossel. Recursive reconstruction on periodic trees. *Random Struct. Algor.*, 13(1):81–97, 1998.
- [Mos01] E. Mossel. Reconstruction on trees: beating the second eigenvalue. *Ann. Appl. Probab.*, 11(1):285–300, 2001.
- [Mos03] E. Mossel. On the impossibility of reconstructing ancestral data and phylogenies. *J. Comput. Biol.*, 10(5):669–678, 2003.
- [Mos04] E. Mossel. Phase transitions in phylogeny. *Trans. Amer. Math. Soc.*, 356(6):2379–2404, 2004.
- [MP03] E. Mossel and Y. Peres. Information flow on trees. Ann. Appl. Probab., 13(3):817–844, 2003.
- [MSW04] F. Martinelli, A. Sinclair, and D. Weitz. Glauber dynamics on trees: boundary conditions and mixing time. *Comm. Math. Phys.*, 250(2):301–334, 2004.
- [Nav01] G. Navarro. A guided tour to approximate string matching. *ACM Computing Surveys (CSUR)*, 33(1):31–88, 2001.
- [NBYST] G. Navarro, R. Baeza-Yates, E. Sutinen, and J. Tarhio. Indexing methods for approximate string matching. *Bulletin of the Technical Committee on*, page 19.
- [Ney71] J. Neyman. Molecular studies of evolution: a source of novel statistical problems. In S. S. Gupta and J. Yackel, editors, *Statistical desicion theory and related topics*, pages 1–27. Academic Press, New York, 1971.
- [NHH00] C. Notredame, D.G. Higgins, and J. Heringa. T-coffee: A novel method for fast and accurate multiple sequence alignment. 2000.
- [RE08] Elena Rivas and Sean R. Eddy. Probabilistic phylogenetic inference with insertions and deletions. *PLoS Comput Biol*, 4(9):e1000172, 09 2008.
- [Roc07] S. Roch. *Markov Models on Trees: Reconstruction and Applications*. PhD thesis, UC Berkeley, 2007.
- [Roc08] Sébastien Roch. Sequence-length requirement for distance-based phylogeny reconstruction: Breaking the polynomial barrier. In *FOCS*, pages 729–738, 2008.
- [Sly09a] A. Sly. Reconstruction of random colourings. *Communications in Mathematical Physics*, 288(3):943–961, 2009.
- [Sly09b] A. Sly. Spatial and Temporal Mixing of Gibbs Measures. PhD thesis, UC Berkeley, 2009.
- [Sly09c] Allan Sly. Reconstruction for the Potts model. In *STOC '09: Proceedings of the 41st annual ACM symposium on Theory of computing*, pages 581–590, New York, NY, USA, 2009. ACM.
- [SR06] Marc A. Suchard and Benjamin D. Redelings. BAli-Phy: simultaneous Bayesian inference of alignment and phylogeny. *Bioinformatics*, 22(16):2047–2048, 2006.
- [SS03] C. Semple and M. Steel. *Phylogenetics*, volume 22 of *Mathematics and its Applications series*. Oxford University Press, 2003.

- [Tho04] J. W. Thornton. Resurrecting ancient genes: experimental analysis of extinct molecules. *Nat. Rev. Genet.*, 5(5):366–375, 2004.
- [TKF91] Jeffrey L. Thorne, Hirohisa Kishino, and Joseph Felsenstein. An evolutionary model for maximum likelihood alignment of dna sequences. *Journal of Molecular Evolution*, 33(2):114–124, 1991.
- [TKF92] Jeffrey L. Thorne, Hirohisa Kishino, and Joseph Felsenstein. Inching toward reality: An improved likelihood model of sequence evolution. *Journal of Molecular Evolution*, 34(1):3–16, 1992.
- [VS08] Krishnamurthy Viswanathan and Ram Swaminathan. Improved string reconstruction over insertiondeletion channels. In SODA '08: Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms, pages 399–408, Philadelphia, PA, USA, 2008. Society for Industrial and Applied Mathematics.
- [WJ94] Lusheng Wang and Tao Jiang. On the complexity of multiple sequence alignment. *Journal of Computational Biology*, 1(4):337–348, 1994.
- [WSH08] Karen M. Wong, Marc A. Suchard, and John P. Huelsenbeck. Alignment Uncertainty and Genomic Analysis. *Science*, 319(5862):473–476, 2008.

A Algorithm

1. **Input.** Children sequences $\hat{x}^1, \ldots, \hat{x}^d$. 2. Initialization. Set $\hat{s}_i(0) := 0, \forall i, \ell = k^{1/3}, r = 1, \text{ and } t = \ell$. 3. Main loop. While $\hat{x}^i[t+\hat{s}_i(r-1)+1:t+\hat{s}_i(r-1)+a]$ is non-empty for all *i*, (a) **Current position.** Set $t = \ell r$. (b) Anchor definition. For each *i*, set $\widehat{A}_r^i = \hat{x}^i [t + \hat{s}_i(r-1) + 1 : t + \hat{s}_i(r-1) + a]$. We say that \widehat{A}_{r}^{i} is the r'th anchor of the i'th child. (If the remaining sequences are not long enough to produce an anchor of length a, we repeat the previous step with the full remaining sequences.) (c) Alignment. For each anchor, we define the set of anchors which agree with it. Formally, $G_r^i = \{ j \in [d], \operatorname{Corr}(\widehat{A}_r^i, \widehat{A}_r^j) > \gamma \}.$ (d) Update. We define the set of aligned children $G_r = \{i : |G_r^i| \ge d-2\}.$ i. Aligned anchors. For each $i \in G_r$, set $\hat{s}_i(r) = \hat{s}_i(r-1)$. ii. Misaligned anchors. For each $i \notin G_r$ define two strings $\widehat{D}_r^i = \hat{x}^i [t + \hat{s}_i (r-1)]$: $t + \hat{s}_i(r-1) + a - 1$ and $\widehat{I}_r^i = \hat{x}^i[t + \hat{s}_i(r-1) + 2:t + \hat{s}_i(r-1) + a + 1]$. If $|\{j \in [d] - \{i\} : \operatorname{Corr}(\widehat{D}_r^i, \widehat{A}_r^j) \ge \gamma\}| \ge d - 2,$ set $\hat{s}_i(r) = \hat{s}_i(r-1) - 1$. If $|\{j \in [d] - \{i\} : \operatorname{Corr}(\widehat{I}_r^i, \widehat{A}_r^j) \ge \gamma\}| \ge d - 2,$ set $\hat{s}_i(r) = \hat{s}_i(r-1) + 1$. (e) Ancestral sequence. Compute $\hat{x}_{t-\ell+1}^0, \dots \hat{x}_t^0$ by performing a sitewise majority on the children in G_r . (If the remaining children sequences are too short to produce a full island, we use whatever is left which should all have equal length by our proof.) (f) **Increment.** Set r := r + 1. 4. **Output.** Output \hat{x}^0 and set \hat{k}_0 to its length.

Figure 1: This is the basic recursive step of our reconstruction algorithm. It takes as input the d inferred sequences of the children $\hat{x}^1, \ldots, \hat{x}^d$ and computes a sequence for the parent \hat{x}^0 . If any of the steps above cannot be accomplished, we abort the reconstruction of the parent and declare it radioactive.

B Further Lemmas

For α going to 0, we have more precisely:

Lemma B.1 (Limit $\alpha \to 0$) Condition on \mathcal{L} . Let

$$\alpha = \frac{1}{h(n)},$$

for $h(n) = \omega(1)$. Then, for n large enough, the root is the father of a (d-1)-ary stable subtree with probability at least

$$1 - \chi = 1 - \frac{1}{\sqrt{h(n)}}.$$

Proof of Lemma B.1: Plugging $\alpha = 1/h(n)$ and $\nu = 1 - 1/\sqrt{h(n)}$ into the recursion derived in the proof of Lemma 3.4, we get

$$(1-\alpha)g(\nu) = \left(1-\frac{1}{h(n)}\right)\left(1-\frac{d}{\sqrt{h(n)}}+O\left(\frac{1}{h(n)}\right)\right)$$
$$+d\left(1-\frac{d-1}{\sqrt{h(n)}}+O\left(\frac{1}{h(n)}\right)\right)\frac{1}{\sqrt{h(n)}}\right)$$
$$= \left(1-\frac{1}{h(n)}\right)\left(1-O\left(\frac{1}{h(n)}\right)\right)$$
$$\geq 1-1/\sqrt{h(n)},$$

for $n \to +\infty$.

C Proofs

Proof of Lemma 3.1: We prove the upper bound by assuming there is no deletion. The lower bound can be proved similarly. The proof goes by induction. Let v be a node at graph distance i from the root. We show that there is C'' > 0 independent of i such that

$$k_v \le k + i\sqrt{C''k\log n}.$$

Since the depth of T is $O(\log n)$, this implies the main claim as long as

$$\sqrt{C''k\log n}\log n \le \zeta k,$$

which follows from our assumption for C' > 0 large enough.

The base case of the induction is satisfied trivially. Assume the induction claim holds for v, the parent of u. It suffices to show that the number of new insertions is at most $\sqrt{C''k\log n}$. By our induction hypothesis, the number of insertions is bounded above by a binomial Z with parameters $k + (i-1)\sqrt{C''k\log n} \le (1+\zeta)k$ and p_{id} w.h.p. By Hoeffding's inequality, taking

$$\eta = \sqrt{\frac{C''' \log n}{(1+\zeta)k}},$$

we have

$$\mathbb{P}[Z > (1+\zeta)kp_{\rm id} + (1+\zeta)k\eta] < \exp(-2((1+\zeta)k\eta)^2/[(1+\zeta)k]) \\ = 1/{\rm poly}(n).$$

By our assumption on p_{id} , we have

$$(1+\zeta)kp_{\rm id} = O\left(\frac{\alpha k^{1/3}}{\log n}\right)$$

so that choosing C'' large enough gives

$$(1+\zeta)kp_{\rm id} + (1+\zeta)k\eta \le \sqrt{C''k\log n}$$

This proves the claim. \blacksquare

Proof of Lemma 3.3: According to Lemma 3.1, the length of the sequence at v is in $[\underline{k}, \overline{k}]$ w.h.p. We denote that event by \mathcal{L}_v . We bound the probability of events $\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_3$ separately.

Let $\overline{N} = \overline{k}/\ell = (1+\zeta)k^{2/3}$. Conditioned on \mathcal{L}_v , there are at most \overline{N} anchors, each of length *a*. By a union bound, the probability that at least one of the sites in the anchors has an indel operation in any child is upper bounded by

$$\begin{aligned} \mathbb{P}[\mathcal{B}_1] &= \mathbb{P}[\mathcal{B}_1 \mid \mathcal{L}_v] \mathbb{P}[\mathcal{L}_v] + \mathbb{P}[\mathcal{B}_1 \mid \mathcal{L}_v^c] \mathbb{P}[\mathcal{L}_v^c] \\ &\leq \overline{N} a dp_{\mathrm{id}} + 1/\mathrm{poly}(n) \\ &= \frac{\alpha a d\overline{N}}{4k^{2/3} a d} + 1/\mathrm{poly}(n) \\ &= \frac{(1+\zeta)k^{2/3}}{k^{2/3}} \cdot \frac{\alpha}{4} + 1/\mathrm{poly}(n) \\ &< \alpha(1/3 - 1/\mathrm{poly}(n)), \end{aligned}$$

where we choose ζ small enough. The quantity we want to estimate is in fact $\mathbb{P}[\mathcal{B}_1 | \mathcal{L}]$ (which is not the same as conditioning on \mathcal{L}_v only). But notice that

$$\mathbb{P}[\mathcal{B}_1] = \mathbb{P}[\mathcal{B}_1 \,|\, \mathcal{L}]\mathbb{P}[\mathcal{L}] + \mathbb{P}[\mathcal{B}_1 \,|\, \mathcal{L}^c]\mathbb{P}[\mathcal{L}^c] \geq \mathbb{P}[\mathcal{B}_1 \,|\, \mathcal{L}]\mathbb{P}[\mathcal{L}],$$

which implies

$$\mathbb{P}[\mathcal{B}_1 \,|\, \mathcal{L}] \le \frac{\alpha(1/3 - 1/\text{poly}(n))}{1 - 1/\text{poly}(n)} < \alpha/3.$$

(This argument shows that it suffices to condition on \mathcal{L}_v . We apply the same trick below.)

To bound the probability of the second event, consider an island I and a son u. The probability that there is an indel when evolving from v to u is at most

$$p_{\rm id}\ell = \frac{\alpha}{4k^{2/3}ad}k^{1/3} = \frac{\alpha}{4k^{1/3}ad}$$

Thus, the probability that more than one child of v experiences an indel in I is at most

$$\begin{split} \sum_{i=2}^{d} \binom{d}{i} \left(\frac{\alpha}{4k^{1/3}ad}\right)^{i} &\leq \sum_{i=2}^{d} \frac{d^{i}}{i!} \left(\frac{\alpha}{4k^{1/3}ad}\right)^{i} \\ &\leq \sum_{i=2}^{d} \frac{1}{i!} \left(\frac{\alpha}{4k^{1/3}a}\right)^{i} \\ &\leq e \left(\frac{\alpha}{4k^{1/3}a}\right)^{2} \\ &= \frac{e\alpha^{2}}{16k^{2/3}a^{2}}, \end{split}$$

where we used that the expression in parenthesis on the second line is < 1. Taking a union bound over all islands, the probability that at least two children experience an indel in the same island is at most

$$\mathbb{P}[\mathcal{B}_2 \,|\, \mathcal{L}] \leq \overline{N} \cdot \frac{e\alpha^2}{16k^{2/3}a^2} \\ = \frac{(1+\zeta)e\alpha^2}{16a^2} \\ < \frac{\alpha}{3},$$

where we used that $\alpha < 1$.

For the third event, consider again an island I and a child u. The probability for at least two indel operations in I when evolving from v to u is at most

$$\begin{split} \sum_{i=2}^{2\ell} \binom{2\ell}{i} \left(\frac{\alpha}{4adk^{2/3}}\right)^i &\leq \sum_{i=2}^{2\ell} \frac{1}{i!} \left(\frac{2\ell\alpha}{4adk^{2/3}}\right)^i \\ &\leq \sum_{i=2}^{2\ell} \frac{1}{i!} \left(\frac{\alpha}{2adk^{1/3}}\right)^i \\ &\leq e \left(\frac{\alpha}{2adk^{1/3}}\right)^2 \\ &\leq \frac{e\alpha^2}{4a^2d^2k^{2/3}}. \end{split}$$

(We use 2ℓ to account for insertions *and* deletions.) Taking a union bound over all islands and children, the probability that there are two indel operations in the same child in the same island is bounded by

$$\mathbb{P}[\mathcal{B}_3 | \mathcal{L}] \leq d\overline{N} \frac{e\alpha^2}{4a^2 d^2 k^{2/3}} \\ \leq \frac{(1+\zeta)e\alpha^2}{4a^2 d} \\ < \alpha/3.$$

Taking a union bound over the three ways in which a site can become radioactive proves the lemma.

Proof of Lemma 3.4: We follow a proof of [Mos01]. Let v be a node at distance r from the leaves. We let ν_r be the probability that v is the root of a (d-1)-ary stable subtree. Let

$$g(\nu) = \nu^d + d\nu^{d-1}(1-\nu).$$

Then, from Lemma 3.3,

$$\nu_r \ge (1 - \alpha)g(\nu_{r-1}).$$

Note that

$$g'(\nu) = d(d-2)\nu^{d-2}(1-\nu).$$

In particular, g is monotone, g(1) = 1, and g'(1) = 0. Hence, there is $1 - \chi < \nu^* < 1$ such that

$$g(\nu^*) > \nu^*.$$

Then, taking

$$1 - \alpha > \nu^* / g(\nu^*),$$

we have

$$\nu_r \ge (1-\alpha)g(\nu_{r-1}) \ge \frac{\nu^*}{g(\nu^*)}g(\nu_{r-1}) \ge \nu^* > 1-\chi,$$

by the induction hypothesis that $\nu_{r-1} \ge \nu^*$. Note in particular that $\nu_0 = 1 \ge \nu^*$.

Proof of Lemma 4.3: Recall that we assume the root state is 0 and all adversarial nodes are 1. Because of the bias towards 1, we cannot apply standard results about recursive majority for symmetric channels [Mos98, Mos04]. Instead, we perform a tailored analysis of this particular channel.

We take asymptotics as $d \to +\infty$ and we show that the probability of reconstruction can be taken to be

$$1-\beta = 1 - \frac{1}{d},$$

for C'' large enough. Let v be the root of $T_{H_0}^{(d)}$. We denote by Z_v the number of non-adversarial children of v in state 0 and by Z'_v the number of nodes among them that return 0 upon applying recursive majority to their respective subtree. Let $q_{H_0}^0$ be the probability of incorrect reconstruction at v (given that the state at v is 0). Then

$$1 - q_{H_0}^0 \geq \mathbb{P}\left[Z'_v \geq \frac{d+1}{2}\right]$$

$$\geq \sum_{i=0}^{d-2} \mathbb{P}\left[Z'_v \geq \frac{d+1}{2} \mid Z_v = i\right] \mathbb{P}[Z_v = i], \qquad (1)$$

where we simply ignored the contribution of the children who flipped to 1.

We prove $q_{H_0}^0 \leq 1/d$ by induction on the height. Let u be a non-adversarial node in $T_{H_0}^{(d)}$ at height h from the leaves to which we associate as above the variables Z_u, Z'_u and the quantity q_h^0 . Note that $q_0^0 = 0$. We assume the induction hypothesis holds for h - 1. Note that conditioned on the state at u being $0 Z_u$ is $Bin(d - 2, (1 - p_s))$ where

$$1 - p_{\rm s} = \frac{1 + \theta_{\rm s}}{2} = \frac{1}{2} + \Theta\left(\sqrt{\frac{\log d}{d}}\right)$$

as $d \to +\infty$. Similarly, given $Z_u = i$, the variable Z'_u is $Bin(i, 1 - q^0_{h-1})$. In particular, the quantity

$$\mathbb{P}\left[Z'_u \ge \frac{d+1}{2} \,|\, Z_u = i\right],\,$$

is monotone in i. We use Chernoff's bound on Z'_u to truncate the lower bound (1). Indeed, let

$$\mu = (1 - p_{\rm s})(d - 2) = \frac{d}{2} + \Upsilon(d),$$

with

$$\Upsilon(d) = \Theta(\sqrt{d\log d}),$$

and

$$\mu(1-\eta) = \frac{d}{2} + \frac{\Upsilon(d)}{2}$$

where in particular

$$\eta = \Theta\left(\sqrt{\frac{\log d}{d}}\right).$$

Then, we have

$$\mathbb{P}[Z_u < \mu(1-\eta)] < \exp(-\mu\eta^2/2) = d^{-\Omega(1)},$$

for C'' large enough. Applying to (1) leads to the lower bound

$$1 - q_h^0 \ge (1 - d^{-\Omega(1)}) \mathbb{P}\left[\operatorname{Bin}\left(\frac{d}{2} + \frac{\Upsilon(d)}{2}, 1 - q_{h-1}^0\right) \ge \frac{d+1}{2} \right].$$

By the induction hypothesis, $q_{h-1}^0 \leq 1/d$. By applying Chernoff's bound again we get

$$\mathbb{P}\left[\operatorname{Bin}\left(\frac{d}{2} + \frac{\Upsilon(d)}{2}, 1 - q_{h-1}^{0}\right) \ge \frac{d+1}{2}\right] > 1 - d^{-\Omega(1)},$$

and therefore $q_h^0 \leq 1/d$. This proves the claim.

Proof of Lemma 4.5: As we pointed out earlier, although the subtrees $(T_{t'}^{**})_{t'=t+1}^{t+m}$ are correlated by the construction of the islands, they are independent of the substitution process. By forcing (randomly) the subtrees $(T_{t'}^{**})_{t'=t+1}^{t+m}$ to be (d-2)-ary and fixing the adversarial nodes to 1, we restore the i.i.d. nature of the sites, from which the result follows.

Proof of Lemma 4.7: This follows from Lemma 4.5, the independence of Λ and Θ , and three applications of Hoeffding's lemma.

Proof of Lemma 4.8: Note that

$$\operatorname{Corr}(\widehat{X},\widehat{Y}) = \frac{1}{m} \sum_{i=1}^{m} \langle \widehat{x}_i \rangle \langle \widehat{y}_i \rangle = \frac{1}{m} \sum_{i=1}^{m} \langle x_i \rangle \langle y_i \rangle \langle z_i \rangle \langle w_i \rangle.$$

Hence,

$$|\operatorname{Corr}(X,Y) - \operatorname{Corr}(\widehat{X},\widehat{Y})| \le \frac{1}{m} \sum_{i=1}^{m} (1 - \langle z_i \rangle \langle w_i \rangle) = 1 - \frac{1}{m} \sum_{i=1}^{m} \langle z_i \rangle \langle w_i \rangle.$$

Now notice by case analysis that

$$\langle z_i \rangle \langle w_i \rangle \ge \langle \lambda_i \rangle \langle \theta_i \rangle - \mathbb{1}_{\langle \lambda_i \rangle = -1} - \mathbb{1}_{\langle \theta_i \rangle = -1}$$

This proves the first claim. The second claim follows from the bounds in Lemma 4.7.

Proof of Lemma 5.1: For the first claim, note that

$$\mathbb{E}[\operatorname{Corr}(\mathscr{A}_r^i, \mathscr{A}_r^j)] = \theta_{\mathrm{s}}^2 = (1 - 2p_{\mathrm{s}})^2$$

where we used that 1) there is no indel in the sites [t + 1 : t + a] between v and i, j; 2) that the sites are perfectly aligned; and 3) that the substitution process is independent of the indel process. We also used the well-known fact that the θ_s 's are multiplicative along a path under our model of substitution [SS03]. The result then follows from Hoeffding's inequality.

For the second claim, because the anchors are now misaligned the t'-th term in $\operatorname{Corr}(\mathscr{A}_r^i, \mathscr{D}_r^j)$ for $t' \in [t+1: t+a]$ is the variable $\langle x_{t'+s_i(r)}^i \rangle \langle x_{t'+s_j(r)-1}^j \rangle$ which is uniform in $\{-1,+1\}$. In particular, we now have

$$\mathbb{E}[\operatorname{Corr}(\mathscr{A}_r^i, \mathscr{D}_r^j)] = 0.$$

The result follows from the method of bounded differences with respect to the independent vectors

$$\{(x_{t'+s_i(r)}^i, x_{t'+s_i(r)}^j)\}_{t'=t}^{t+a}.$$

Proof of Lemma 5.2: This follows from Lemmas 4.8 and 5.1 and the triangle inequality. ■

D Completing the Proof of the Main Theorem

Having conditioned on the event \mathcal{E} , we justify the correctness of our reconstruction method via the following induction. The top level of the induction establishes Theorem 1.3.

Induction hypothesis. Consider a parent v in T^* ; in particular, v is stable. We assume that the following conditions, denoted by (\star) , are satisfied: For all children $i \in [d]$ of v belonging to T^*

1. Alignment. For all children i' of i with $i' \in T^*$ and all $r = 1, \ldots, \bar{k}/\ell - 1$,

$$\hat{s}_{i'}(r) = s_{i'}(r).$$
 (2)

(This condition is trivially satisfied for values of $r\ell$ that are larger than the sequence length of i'.)

2. **Reconstruction.** Moreover, we have $\hat{k}_i = k_i$ and for all $t = 1, ..., k_i$, the following holds:

Let L_i be the leaves below i with $n_i = |L_i|$. Let H be the level of v. Let L_t^{**} be the gateway leaves for site t. For $u \in L_t^{**}$ let $F_u(t)$ be the position of site t in u. Note that \hat{x}_t^i can be written as $\hat{x}_t^i = \text{Maj}^{H-1}(z_1, \ldots, z_{n_i})$, where z_j is either \sharp or $x_{\flat_j}^j$ for an appropriate function \flat_j . Our hypothesis is that

$$\forall u \in L_t^{**}, \ b_u = F_u(t). \tag{3}$$

In particular, the ancestral reconstruction \hat{X}_i dominates the adversarial reconstruction \hat{X}'_i .

The base case where v is a leaf is trivially satisfied.

Alignment. We begin with the correctness of the alignment.

Lemma D.1 (Induction: Alignment) Assuming \mathcal{E} and (\star) , the algorithm infers s_i correctly for all children $i \in [d]$ which are also in T^* , that is, (2) holds for v.

Proof of Lemma D.1: Let Π denote the set of children of v in T^* . The proof follows by induction on r. The base case r = 0 is trivial. Assume correctness for r - 1.

If there is no indel in any of the children $i \in \Pi$ between the sites $(r-1)\ell$ and $r\ell$ of v, then under \mathcal{E} , (\star) and Lemma 5.2 we have $\Pi \subseteq G_r$. In that case, for all $i \in \Pi$ we have $\hat{s}_i(r) = \hat{s}_i(r-1) = s_i(r-1) = s_i(r)$, where the second equality is from (\star) .

If there is an indel operation in island r, then since v is stable only one indel operation occurred in one child. Denote the child with an indel by j. Assume the indel is a deletion. (The case of the insertion is handled similarly.) If j is not in T^* we are back to the previous case. So assume j is in T^* . Again, from \mathcal{E} , (\star) and Lemma 5.2 the other children in T^* are added to the set G_r , and the shift value will be computed correctly for them. Moreover by (\star), for every $i \in \Pi - \{j\}$,

$$f_i(r\ell + 1) = r\ell + 1 + s_i(r) = r\ell + 1 + \hat{s}_i(r) = r\ell + 1 + \hat{s}_i(r-1)$$

which is the starting point of \widehat{A}_r^i . Also,

$$f_j(r\ell + 1) = r\ell + 1 + s_j(r) = r\ell + 1 + s_j(r-1) - 1 = r\ell + 1 + \hat{s}_j(r-1) - 1 = r\ell + \hat{s}_j(r-1),$$

which is the starting point of \widehat{D}_r^j . Thus according to Lemma 5.2 \widehat{D}_r^j matches \widehat{A}_r^i for all $i \in \Pi \cap G_r$. As there are d-2 children in $\Pi \cap G_r$, we get that the algorithm sets

$$\hat{s}_j(r) = \hat{s}_j(r-1) - 1 = s_j(r-1) - 1 = s_j(r),$$

as required. Note also that in this case, according to Lemma 5.2 again, \widehat{A}_r^j does not have high correlation with \widehat{A}_r^i for any $i \in \Pi \cap G_r$, and thus we will consider \widehat{I}_r^j and \widehat{D}_r^j . Similarly, \widehat{I}_r^j does not have high correlation with \widehat{A}_r^i for any $i \in \Pi \cap G_r$, and thus we will not try to set $\widehat{s}_i(r)$ twice.

Ancestral reconstruction. We use Lemma D.1 to prove that the ancestral reconstruction dominates the adversarial reconstruction. In the algorithm, we perform a sitewise majority vote over the children of v in G_r (these are the aligned children—see the description of the algorithm in Figure 1). For notational convenience, we assume that in fact we perform a majority vote over *all* children but we replace the states of the children outside G_r with \sharp .

Lemma D.2 (Induction: Reconstruction) Assuming \mathcal{E} , (*) and the conclusion of Lemma D.1, (3) holds for v. In particular, the ancestral reconstruction \widehat{X}_v dominates the adversarial reconstruction \widehat{X}'_v .

Proof of Lemma D.2: The second claim follows from the first one together with the construction of the adversarial process and the monotonicity of majority.

As for the first claim, by Lemma D.1 for each site of v there are d-2 uncorrupted children islands containing this site such that the children are also in T^* . In particular, the d-2 corresponding sites in the children are correctly aligned. Moreover, by the induction hypothesis, each corresponding site in the children satisfy (3). By taking a majority vote over these sites we get (3) for v as well.

A small technical detail is handling the case where the last island has less than a sites, and thus does not contain an anchor. However, in this case, if the father is stable then there are no indel operations at all in the last island, and therefore aligning it according to the previous one gives the right result.