

RECONSTRUCTING TREES FROM TRACES

BY SAMI DAVIES¹, MIKLÓS Z. RÁCZ² AND CYRUS RASHTCHIAN³

¹*Department of Mathematics, University of Washington, daviess@uw.edu*

²*Princeton University, mrac@princeton.edu*

³*Department of Computer Science & Engineering, University of California, San Diego, crashtchian@eng.ucsd.edu*

We study the problem of learning a node-labeled tree given independent traces from an appropriately defined deletion channel. This problem, tree trace reconstruction, generalizes string trace reconstruction, which corresponds to the tree being a path. For many classes of trees, including complete trees and spiders, we provide algorithms that reconstruct the labels using only a polynomial number of traces. This exhibits a stark contrast to known results on string trace reconstruction, which require exponentially many traces, and where a central open problem is to determine whether a polynomial number of traces suffice. Our techniques combine novel combinatorial and complex analytic methods.

1. Introduction. Statistical reconstruction problems aim to recover unknown objects given only noisy samples of the data. In the *string trace reconstruction* problem, there is an unknown binary string, and we observe noisy samples of this string after it has gone through a deletion channel. This deletion channel independently deletes each bit with constant probability q and concatenates the remaining bits. The channel preserves bit order, so we observe a sampled subsequence known as a *trace*. The goal is to learn the original string with high probability using as few traces as possible.

The string trace reconstruction problem (with insertions, substitutions, and deletions) directly appears in the problem of DNA data storage [9, 12, 17, 19, 35, 38, 39]. It is crucial to minimize the sample complexity, as this directly impacts the cost of retrieving data stored in synthetic DNA. Since there is an exponential gap between upper and lower bounds for the string trace reconstruction problem, it is motivating to study variants. We introduce a generalization of string trace reconstruction called *tree trace reconstruction*, where the goal is to learn a node-labelled tree given traces from a deletion channel. From a technical point of view, tree trace reconstruction may aid in understanding the interplay of combinatorial and analytic approaches to reconstruction problems and can be a springboard for new ideas. From an applications point of view, current research on DNA nanotechnology has demonstrated that structures of DNA molecules can be constructed into trees and lattices. In fact, recent research has shown how to distinguish different molecular topologies, such as spiders with three arms, from line DNA using nanopores [24]. These results may open the door for other tree structures and be useful for applications like DNA data storage.

Let X be a rooted tree with unknown binary labels on its n nonroot nodes. The goal of tree trace reconstruction is to learn the labels of X with high probability, using the minimum number of traces, knowing only q , the deletion model, and the structure of X .

We consider two deletion models. In both models, each nonroot node v in X is deleted independently with constant probability q —the root is never deleted—and deletions are associative. The resulting tree is called a trace. We assume that X has a canonical ordering of

Received June 2019; revised August 2020.

MSC2020 subject classifications. Primary 60C05, 68Q32; secondary 30C80.

Key words and phrases. Trace reconstruction, tree trace reconstruction, deletion channel, Littlewood polynomials.

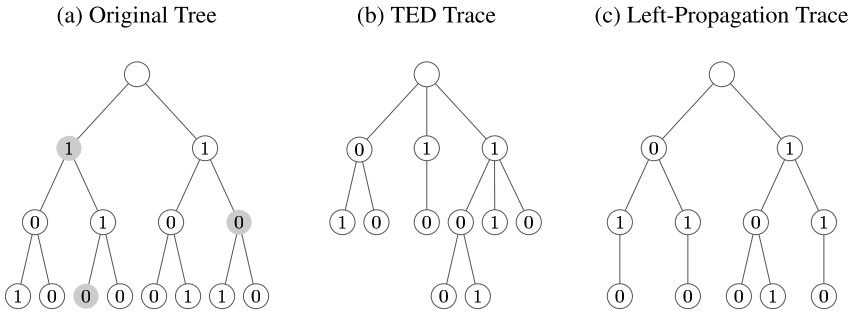


FIG. 1. Deletion models. Gray nodes deleted from original tree (a). Resulting trace in the TED Model (b) and the Left-Propagation Model (c).

its nodes, and the children of a node have a left-to-right ordering. For the left-propagation model, we define the *left-only* path starting at v as the path that recursively goes from parent to left-most child.

- *Tree edit distance (TED) model:* When v is deleted, all children of v become children of v 's parent. Equivalently, contract the edge between v and its parent, retaining the parent's label. The children of v take v 's place as a continuous subsequence in the left-to-right order.
- *Left-propagation model:* When v is deleted, recursively replace every node (together with its label) in the left-only path starting at v with its child in the path. This results in the deletion of the last node of the left-only path, with the remaining tree structure unchanged.¹

Figure 1 depicts traces in both deletion models for a given original tree and set of deleted nodes. When X is a path or a star, then both models coincide with the string deletion channel. After posting this paper to arXiv, subsequent work has shown that these are the most difficult trees to reconstruct in terms of sample complexity [30]. In other words, the sample complexity to reconstruct an arbitrarily labelled tree on n nodes is no more than the sample complexity to reconstruct an arbitrarily labelled string on n bits.

A key motivation for the tree edit distance model is that deletions in the TED model correspond exactly to the deletion operation in tree edit distance, which is a well-studied metric for pairs of labeled trees used in applications [7, 40]. Our main motivation for the left-propagation model is more theoretical: it preserves different structural properties—for instance, a node's number of children does not increase (see Figure 1)—and poses different challenges than the TED model.

1.1. Related work.

Previous results on string trace reconstruction. Introduced by Batu, Kannan, Khanna, and McGregor [6], string trace reconstruction has received a lot of attention, especially recently [10, 14, 15, 20–23, 31, 34, 37]. Yet there is still an exponential gap between the known upper and lower bounds for the number of traces needed to reconstruct an arbitrary string with high probability and constant deletion probability: it is known that $\exp(O(n^{1/3}))$ traces are sufficient [14, 15, 34] and $\tilde{\Omega}(n^{3/2})$ traces are necessary [10, 21]. Determining whether a polynomial number of traces suffice is a challenging open problem in the area. A well-studied variant is reconstructing a string with random, average-case labels, instead of arbitrary, worst-case labels [6, 22]. This is relevant for applications to DNA data storage [35].

¹Since the BFS order on X is arbitrary (but fixed), the choice of using the left-only path (as opposed to, say, the right-only one) does not *a priori* bias certain nodes.

In a few of our algorithms, we will reduce various subproblems to the string trace reconstruction problem, and hence, we will use existing results as a black box. For future reference, we precisely state the previous results now. Let $T(n, \delta)$ and $\widehat{T}(n, \delta)$ denote the minimum number of traces needed to reconstruct an n -bit worst-case and average-case string, respectively, with probability at least $1 - \delta$, where the dependence on the deletion probability q is left implicit.

THEOREM 1.1 ([14, 15, 34]). *The number of traces $T(n, \delta)$ needed to reconstruct a worst-case n -bit string with probability $1 - \delta$ satisfies $T(n, \delta) \leq \ln(\frac{1}{\delta}) \cdot e^{Cn^{1/3}}$, for C depending on q .*

THEOREM 1.2 ([22]). *The number of traces $\widehat{T}(n, \delta)$ needed to reconstruct a random n -bit string with probability $1 - \delta$ satisfies $\widehat{T}(n, \delta) \leq \ln(\frac{1}{\delta}) \cdot e^{C \log^{1/3}(n)}$, for C depending on q .*

In terms of lower bounds, it is known that $T(n, \delta) = \widetilde{\Omega}(n^{1.5})$ and $\widehat{T}(n, \delta) = \widetilde{\Omega}(\log^{5/2}(n))$, for any δ bounded away from one [10, 21]. The proofs of Theorem 1.1 rely on a mean-based algorithm, one only using the mean of single bits from traces, and the bound is optimal for mean-based algorithms [14, 15, 34].

Other variants of trace reconstruction. Due to the exponential gap between upper and lower bounds in the string trace reconstruction problem, an array of variants have been studied recently. Cheraghchi, Gabrys, Milenkovic, and Ribeiro introduce the study of coded trace reconstruction, where the goal is to design efficiently encodable codes whose codewords can be efficiently reconstructed with high probability [11]. Krishnamurthy, Mazumdar, McGregor, and Pal study trace reconstruction on matrices, where rows and columns of a matrix are deleted and a trace is the resulting submatrix. They also study string trace reconstruction on sparse strings [27]. Ancestral state reconstruction is a generalization of string trace reconstruction, where traces are no longer independent, but instead evolve based on a Markov chain [4].

There is also a deterministic version of string trace reconstruction [29]. Let the k -deck of a string be the multiset of its length k subsequences. The question is to establish how large k must be to uniquely determine an arbitrary string of n bits. Currently, the best known bounds stand at $k = O(\sqrt{n})$ and $k = \exp(\Omega(\sqrt{\log n}))$, due respectively to Krasikov and Roditty [26] and Dudík and Schulman [16]. This result has also been used to study population recovery, the problem of learning an unknown distribution of bit strings given noisy samples from the distribution [5].

The term trace complexity has appeared in a network inference context, but the models and definition of a trace are incomparable to ours [1]. Other results on deletion channels appear in the survey by Mitzenmacher [32].

Other graph reconstruction models. While we are unaware of previous work on reconstructing trees using traces (besides strings), a large variety of other graph-centric reconstruction problems have been considered.

The famous *reconstruction conjecture*, due to Kelly [25] and Ulam [36], posits that every graph G is uniquely determined by its *deck*, where the deck of G is the multiset of subgraphs obtained by deleting a single vertex from G . Here, the (sub)graphs are unlabeled, and the goal is to determine G up to isomorphism. The reconstruction conjecture remains open, although it is known for special cases, such as trees and regular graphs [25, 28].

Mossel and Ross introduced and studied the shotgun assembly problem on graphs, where they use small vertex-neighborhoods to uniquely identify an unknown graph [33].

1.2. *Our results.* We provide algorithms for two main classes of trees: complete k -ary trees and spiders. In a *complete k -ary tree*, every nonleaf node has exactly k children, and all leaves have the same depth. An (n, d) -*spider* consists of n/d paths of $d + 1$ nodes, all starting from the same root. Figure 11 depicts an example spider, and it demonstrates that both deletion models lead to the same trace for spiders. We focus on these two classes because of their varying amount of structure. Spiders behave like a union of disjoint paths, except when some paths have all of their nodes deleted. This allows us to extend methods from string trace reconstruction, with a slightly more complicated analysis. On the other hand, complete k -ary trees are so structured that we can use more combinatorial algorithms, which have proven less successful for string trace reconstruction so far. We believe our methods could be used to prove results for larger classes of trees, as well.

In what follows, we use *with high probability* to mean with probability at least $1 - O(1/n)$. Also, we let $[t]$ for $t \in \mathbb{N}$ denote the set $\{1, 2, \dots, t\}$.

1.2.1. *TED model for complete k -ary trees.* Let X be a rooted complete k -ary tree along with unknown binary labels on its n nonroot nodes. Since $k = 1$ and $k = n$ are identical to string trace reconstruction, we focus on $1 < k < n$. We provide two algorithms to reconstruct X , depending on whether the degree k is large or small.

We state our theorems in terms of $T(k, \delta)$, since our reductions use algorithms for string trace reconstruction as a black box and the current bounds on $T(k, \delta)$ may improve in the future.

THEOREM 1.3. *In the TED model, there exist $c, c' > 0$ depending only on q such that if $k \geq c \log^2(n)$, then it is possible to reconstruct a complete k -ary tree on n nodes with $\exp(c' \cdot \log_k n) \cdot T(k, 1/n^2)$ traces with high probability.*

Theorem 1.1 implies that $T(k, 1/n^2) = \exp(O(k^{1/3}))$ if $k \geq c \log^2(n)$, so the trace complexity in Theorem 1.3 is currently $\exp(O(\log_k(n) + k^{1/3}))$. This is $\text{poly}(n)$ as long as $k = O(\log^3 n)$.

THEOREM 1.4. *In the TED model, there exists $C > 0$ depending only on q such that $\exp(Ck \log_k n)$ traces suffice to reconstruct a complete k -ary tree on n nodes with high probability.*

In particular, when k is a constant, then the trace complexity of Theorem 1.4 is $\text{poly}(n)$. Theorem 1.4 makes no restrictions on k , but uses more traces than Theorem 1.3 for $k \geq c \log^2 n$. It would be desirable to smooth out the dependence on k between our two theorems. In particular, we leave it as an intriguing open question to determine whether $\text{poly}(n)$ traces suffice for all $k \leq \log^3(n)$.

1.2.2. *Left-propagation model for complete k -ary trees.* We provide two reconstruction algorithms for k -ary trees in the left-propagation model, leading to the following two theorems.

THEOREM 1.5. *In the left-propagation model, there exists $c > 0$ depending only on q such that if $k \geq c \log n$, then $T(d + k, 1/n^2)$ traces suffice to reconstruct a complete k -ary tree of depth $d = O(\log_k n)$ with high probability.*

When $k \geq c \log n$, then $d + k = O(k)$, and we can reconstruct an n -node complete k -ary tree with $\exp(O(k^{1/3}))$ traces by using Theorem 1.1.

We also provide an alternate algorithm that makes no assumptions on k .

THEOREM 1.6. *In the left-propagation model, $O(n^\gamma \log n)$ traces suffice to reconstruct an n -node complete k -ary tree with high probability, where $\gamma = \ln(\frac{1}{1-q})(\frac{c'k}{\ln n} + \frac{1}{\ln k})$, for a constant $c' > 1$.*

Theorem 1.6 implies that $\text{poly}(n)$ traces suffice to reconstruct a k -ary tree whenever $k = O(\log n)$ and q is a constant. Moreover, for small enough q and k , the algorithm needs only a sublinear number of traces (e.g., binary trees with $q < 1/2 - \varepsilon$). From Theorem 1.1, the bound in Theorem 1.6 can be more simply thought of as $\exp(C' \cdot (d + k))$; and, in Theorem 1.5 as $\exp(C \cdot (d + k)^{1/3})$.

1.2.3. Spiders. Recall that the TED and left-propagation deletion models are the same for spiders. We provide two reconstruction algorithms, depending on whether the depth d is large or small.

THEOREM 1.7. *Assume that $d \leq \log_{1/q} n$. For $q < 0.7$, there exists $C > 0$ depending only on q such that $\exp(C \cdot d(nq^d)^{1/3})$ traces suffice to reconstruct an (n, d) -spider with high probability.*

To understand the statement of this theorem, consider $d = c \log_{1/q} n$ with $c < 1$. A black-box reduction to the string case results in using $\exp(\tilde{\Omega}(n^{1-c}))$ traces for reconstruction (see Section 5.4), whereas Theorem 1.7 improves this to $\exp(\tilde{O}(n^{(1-c)/3}))$.

Theorem 1.7 actually extends to any deletion probability $q < 1/\sqrt{2} \approx 0.707$, but this requires taking d to be larger than some constant depending on q . We discuss further in Remark 3 why the regime of $q > 1/\sqrt{2}$ is difficult to handle. Our approach extends previous results based on complex analysis [14, 15, 34]. As the main technical ingredient, we prove new bounds on certain polynomials whose coefficients are small in modulus. In particular, we analyze a generating function that might be of independent interest, related to Littlewood polynomials.

For large depth $d \geq \log_{1/q} n$, full paths of the spider are unlikely to be completely deleted, and we derive the following result via a reduction to string trace reconstruction.

PROPOSITION 1.8. *For $q < 1$ and all n large enough, an (n, d) -spider with $d \geq \log_{1/q} n$ can be reconstructed with $2 \cdot T(d, \frac{1}{2n^2})$ traces with high probability.*

Using Theorem 1.1, the current bound for Proposition 1.8 is $2 \cdot T(d, \frac{1}{2n^2}) \leq \exp(O(d^{1/3}))$. Comparing Theorem 1.7 and Proposition 1.8, we see that the bounds in the exponent are $d(nq^d)^{1/3}$ and $d^{1/3}$, for $d \leq \log_{1/q} n$ and $d \geq \log_{1/q} n$, respectively. We leave it as an open question to unify these bounds, and in particular, to determine whether the jump is necessary as d crosses $\log_{1/q} n$.

1.2.4. Average-case labels for trees. Our results have focused on trees with worst-case, arbitrary labels. Assuming the binary labels are uniformly distributed independent bits leads to significantly improved bounds. For the string case, Theorem 1.2 implies that $\hat{T}(k, 1/n^2) = \exp(O(\log^{1/3} k + \log \log n))$ traces suffice to reconstruct a random binary string with high probability. For three of our results, we can use this as a black box and replace the dependence on $T(k, 1/n^2)$ with $\hat{T}(k, 1/n^2)$ for average-case labeled trees. The average-case trace complexity for k -ary trees under the TED model— analogously to Theorem 1.3— becomes $\exp(O(\log_k(n) + \log^{1/3} k))$ when $k \geq c \log^2(n)$. For the left-propagation model— analogously to Theorem 1.5—the average-case trace complexity becomes $\exp(O(\log^{1/3} k +$

$\log \log n$) when $k \geq c \log n$. For (n, d) -spiders with depth $d \geq \log_{1/q} n$ —analogously to Proposition 1.8—the average-case trace complexity becomes $\exp(O(\log^{1/3} d + \log \log n))$. Since it is straightforward to use the average-case string result instead of the worst-case result to obtain the results just described, we restrict our exposition to worst-case labeled k -ary trees and spiders.

1.3. *Overview of TED deletion algorithms.* Previous work on string trace reconstruction mostly utilizes two classes of algorithms: mean-based methods, which use single-bit statistics for each position in the trace, and alignment-based methods, which attempt to reposition subsequences in the traces to their true positions.

Although mean-based algorithms are currently quantitatively better for string reconstruction, they seem difficult to extend to k -ary trees under the TED deletion model. Specifically, mean-based methods require a precise understanding of how the bit in position j' of the original tree affects the bit in position j of the trace. For strings, there is a global ordering of the nodes which enables this. Unfortunately, for k -ary trees with $k \notin \{1, n\}$ under the TED model, nodes may shift to a variety of locations, making it unclear how to characterize bit-wise statistics. To circumvent this challenge, we provide two new algorithms, depending on whether or not the degree k is large ($k \geq c \log^2(n)$). The main idea is to partition the original tree into small subtrees and learn their labels using a number of traces parameterized primarily by k and $\log_k n$, which can be much smaller than n .

When k is large enough, we will be able to localize root-to-leaf paths, in the sense that we can identify the location of their nonleaf nodes in the original tree with high probability. By covering the internal nodes of the tree by such paths, we will directly learn the labels for all nonleaf nodes. Then, we observe that the leaves can be naturally partitioned into stars of size k , and we can learn their labels by reducing to string trace reconstruction (for strings on k bits). Any improvement to string trace reconstruction will lead to a direct improvement for k -ary trees with large degree.

When k is small, our localization method fails, and we resort to looking at traces which contain even more structure (which requires more traces). We decompose the entire tree into certain subtrees and recover their labels separately. We define a property which is easily detectable among traces and show that when this property holds, we can extract labels for the subtrees that are correct with probability at least $2/3$. Then, we take a majority vote to get the correct labels with high probability.

1.4. *Overview of left-propagation algorithms.* As with the TED model, we combine mean-based and alignment-based strategies, and we provide different algorithms depending on whether the degree is large or small. The two algorithms differ in how they align certain subtrees of traces to positions in the tree.

When k is large enough ($k \geq c \log n$ for a constant $c > 0$), our first algorithm will use results from string trace reconstruction as a black box. The key idea is that certain subtrees will behave as if they were strings on $O(k)$ bits in the string deletion model. Although this does not happen in all traces, we show that it occurs with high probability. Overall, we partition X into such subtrees, and we reduce to string reconstruction results to recover the labels separately.

On the other hand, when k is small (such as binary trees with $k = 2$), we do not know how to reduce to string reconstruction. Instead, our second algorithm waits until a larger subtree survives in a trace. We show that this makes the alignment essentially trivial, and we can directly recover the labels for certain subtrees. Quantitatively, the trace complexity of the first algorithm is better, but the reconstruction only succeeds for large enough k .

1.5. *Overview of spider techniques.* When the paths of a spider are sufficiently long—specifically, if they have depth $d \geq \log_{1/q} n$ —then with probability close to 1, no path is fully deleted in a given trace. This allows us to trivially match paths of the trace spider to paths of the original spider and then use string trace reconstruction algorithms on the individual paths, leading to Proposition 1.8.

When the paths of a spider are shorter ($d < \log_{1/q} n$), many traces have paths fully deleted. As illustrated in Figure 11, when paths are fully deleted from a spider, it is unclear which paths were deleted, which forces us to align paths from different traces. We bypass direct alignment-based methods and instead use a mean-based algorithm that generalizes the methods introduced in the proof of Theorem 1.1 by [14, 15, 34]. The main difficulty we address is that, in contrast to strings which are one-dimensional, spiders are two-dimensional: one dimension representing which path in the spider a node is in, and the other representing where in a path a node is.

1.6. *Outline.* The rest of the paper is organized as follows. Preliminaries are in Section 2. The proofs of Theorem 1.3 and Theorem 1.4 for k -ary trees under the TED model appear in Section 3. The proofs of Theorem 1.5 and Theorem 1.6 for the Left-Propagation model appear in Section 4. The spider reconstruction preliminaries and algorithms for Theorem 1.7 and Proposition 1.8 are in Section 5. The three main sections can be read independently, after their preliminaries. We conclude in Section 6.

2. Preliminaries. In what follows, X denotes the (known) underlying tree, along with the (unknown) binary labels on its n nonroot nodes.

Standard tree definitions. We say that X is *rooted* if it has a fixed root node. We assume the root is never deleted (for further explanation see Remark 4). An *ancestor* (resp. *descendant*) of a node v is a node reachable from v by proceeding repeatedly from child to parent (resp. parent to child). We say v is a *leaf* if it has no children, and otherwise v is an *internal* node. The *length* of a path equals the number of nodes in it. The *depth* of v is the number of edges in the path from the root to v . The *height* of v is the number of edges in the longest path between v and a leaf. The depth of a rooted tree is the height of the root. We say that X is a complete k -ary tree of depth d if every internal node has k children and all leaves have depth d .

2.1. *k -ary tree algorithm preliminaries.* Let X be a rooted complete k -ary tree with depth d . We index the nonroot nodes according to the BFS order on X (the root is not indexed; the children of the root are $\{0, 1, \dots, k-1\}$, etc.). We identify nodes of X with their index. For $t \in [d]$, let \mathcal{J}_t be the nodes at depth t . Define $\mathcal{I}_1 := \mathcal{J}_1 = \{0, 1, \dots, k-1\}$, and for $t \geq 2$,

$$\mathcal{I}_t := \{i \in \mathcal{J}_t \mid i \bmod k \neq 0\}.$$

In words, for $t \geq 2$, \mathcal{I}_t is the set of nodes at depth t which are not left-most among their siblings. Define also $\mathcal{I} := \bigcup_{t=1}^{d-1} \mathcal{I}_t$.

We define three unlabeled subtrees of X . Let $P_X(i)$ be the path from the root to i in X . Define $H_X(i)$ as the union of the left-only path starting at i , descending to a leaf ℓ , and the $k-1$ siblings of ℓ . Finally, define $G_X(i) := P_X(i) \cup H_X(i)$. See Figure 2 for an example of these subtrees. For clarity, we note that if i has depth t in X (i.e., $i \in \mathcal{J}_t$), then $|P_X(i)| = t+1$ and $|H_X(i)| = d-t+k$ and $|G_X(i)| = d+k$.

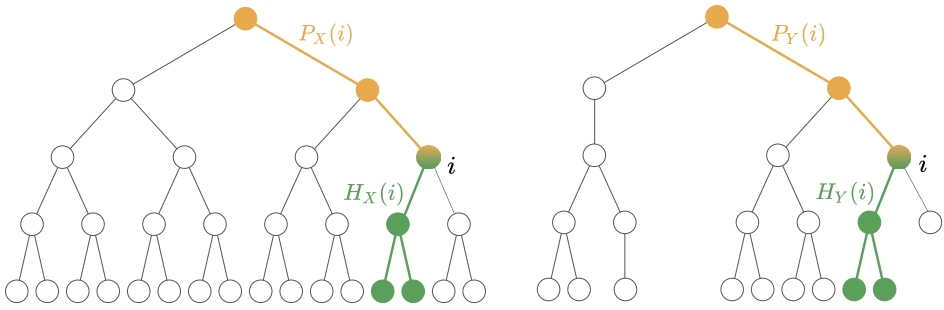


FIG. 2. Canonical subtrees for k -ary trees, in the original tree (left) and trace (right).

Canonical subtrees of traces. We also consider certain subtrees of a trace Y . They will be analogous to $P_X(i)$, $H_X(i)$, and $G_X(i)$, and they only depend on the position of i in X . We will denote them as $P_Y(i)$, $H_Y(i)$, and $G_Y(i)$. Intuitively, they are subtrees in Y obtained by looking at nodes that should be in the same position as the corresponding ones in X . However, the node i does not necessarily belong to these subtrees (e.g., it may have been deleted in Y , or another node may be in its place). In what follows, we refer to subtrees as sequences of nodes in the BFS order, since the edge structure will be clear from context (i.e., the subtree is the induced subgraph on the relevant nodes).

We now formally define $P_Y(i)$, $H_Y(i)$, and $G_Y(i)$, which are also depicted in Figure 2. Fix i , and let u_0, u_1, \dots, u_{d-1} be the internal nodes in $G_X(i)$, where u_t has depth t , and let u_d, \dots, u_{d+k-1} be the leaf nodes, ordered left-to-right in the BFS order. Define $\pi_i : \{0, 1, \dots, d-1\} \rightarrow \{0, 1, \dots, k-1\}$ so that $\pi_i(t)$ is the position of u_{t+1} in X among its siblings (the children of its parent u_t). Note that π_i is independent of the labels of X . Let t_i be the depth of i in X . We define $P_Y(i)$ as the path v_0, v_1, \dots, v_{t_i} in Y obtained from the following process. Set v_0 to be the root. Then, for $t \in [t_i]$, let v_t be the node at depth t in Y that is in position $\pi_i(t-1)$ among the k children of v_{t-1} , where we abort and set $P_Y(i) = \perp$ if v_{t-1} does not have exactly k children. Similarly, let $G_Y(i)$ be the subtree $v_0, v_1, \dots, v_{d+k-1}$, where v_t is defined as follows. Set v_0 to be the root in Y . Then, for $t \in [d-1]$, let v_t be the node at depth t in Y that is in position $\pi_i(t-1)$ among the k children of v_{t-1} , where we abort and set $G_Y(i) = \perp$ if v_{t-1} does not have exactly k children. Finally, set v_d, \dots, v_{d+k-1} to be the k children of v_{d-1} , and again we set $G_Y(i) = \perp$ if v_{d-1} does not have precisely k children. If $G_Y(i) \neq \perp$, then set $H_Y(i) = v_{t_i}, \dots, v_{d+k-1}$, and otherwise, set $H_Y(i) = \perp$. Observe that if $G_Y(i) \neq \perp$, then we have $G_Y(i) = P_Y(i) \cup H_Y(i)$.

We remark that $G_Y(i)$, $H_Y(i)$, and $P_Y(i)$ depend only on π_i and the tree structure of Y , and therefore they do not use any label information from X . We also note that whether these subtrees are set to \perp will be significant, since this implies certain structural properties of traces. If all nodes in $G_X(i)$ survive in a trace Y , then we say that Y contains $G_X(i)$. We write $G_Y(i) = G_X(i)$ if the nodes in these subtrees are exactly the same (by construction, the edges will also be the same). We conclude this section with two remarks that are useful for reconstruction of X .

REMARK 1. If $G_Y(i) = G_X(i)$, then we can reconstruct the labels of $G_X(i)$ bit by bit by copying the label to be that from the corresponding bit in $G_Y(i)$. The same applies for $H_X(i)$ and $P_X(i)$.

REMARK 2. To reconstruct labels of X , one can reconstruct labels of subtrees of X , where the subtrees cover all nodes of X .

3. Reconstructing trees, TED deletion model. In this section we prove our two results for k -ary trees in the TED model.

3.1. *Proof of Theorem 1.3 concerning large degree trees.* Our algorithm utilizes a structure that occurs when $k \geq c \log^2(n)$. Recall that for a node i in X , we think of i 's children as being ordered consecutively, left-to-right, based on the BFS ordering of X .

DEFINITION 3.1. Let Y be a trace of a tree X . We say that Y is b -balanced if, for every internal node i in X , at most b consecutive children of i have been deleted in Y .

CLAIM 1. If X has n nodes, then a trace Y is b -balanced with probability at least $1 - nq^b$.

PROOF. Any set of b consecutive nodes is deleted with probability q^b . Since there are at most n starting nodes for a run of b nodes, a union bound proves the claim. \square

Since $T(k, 1/n^2) = \exp(O(k^{1/3}))$ by Theorem 1.1, the number of traces used in Theorem 1.3 is $\exp(O(\log_k(n) + k^{1/3}))$. Therefore, setting $b := 10\sqrt{k}/\log(1/q) = \Omega(\log n)$, Claim 1 and a union bound show that with high probability *all* traces will be b -balanced. As we shall see, the benefit of this balanced structure manifests itself in the proof of correctness of the reconstruction algorithm used for Theorem 1.3.

Our reconstruction algorithm that proves Theorem 1.3 consists of two main steps:

1. (*Finding paths and grouping traces*). First, we process all the traces and group them into different sets (which may overlap). This grouping of traces is based on finding root-to-leaf paths that are preserved in a trace and estimating where these paths came from in X . We term this latter algorithm the FINDPATHS algorithm; see Algorithm 1 below for its pseudocode.

2. (*Reconstruction*). We then analyze each subset of traces. Each of these leads to reconstructing the labels for a particular subset of X , consisting of a path from the root to a node at depth $d - 1$, together with the k children of this node. Finally, we output the union of all such labels as the estimated labels of X . In other words, we cover X with the collection of subtrees $\{G_X(j) : j \in \mathcal{J}_{d-1}\}$ and estimate the labels of each $G_X(j)$ separately. Algorithm 2 below states, in pseudocode, the full reconstruction algorithm, which calls Algorithm 1 as a subroutine.

The FINDPATHS algorithm. We start by describing the FINDPATHS algorithm (see Algorithm 1). The input to this algorithm is a trace Y , while the output of the algorithm will be a subset of \mathcal{J}_{d-1} , the nodes of X at depth $d - 1$; the conceptual meaning of this subset will be clear once the algorithm is described. We recall that we index the nonroot nodes according to the BFS order on X , and we will interchangeably refer to nodes and their BFS index.

The first part of the FINDPATHS algorithm is to identify root-to-leaf paths that have been preserved (i.e., no vertex in the path has been deleted) in the trace Y . This is straightforward, since if a root-to-leaf path in X is preserved, then the corresponding leaf has depth d in Y ; and vice versa, every leaf in Y that has depth d corresponds to a root-to-leaf path in X that was preserved. Once all surviving root-to-leaf paths have been identified, we collect in the set \mathcal{S} all the nodes of Y that are on a surviving root-to-leaf path and have depth $d - 1$ (see lines 1–6 of Algorithm 1).

We know that each node $v \in \mathcal{S}$ must have come from a node $w \in \mathcal{J}_{d-1}$ (i.e., a node in X of depth $d - 1$). The second and final part of the FINDPATHS algorithm consists of estimating, for each node $v \in \mathcal{S}$, which original node $w \in \mathcal{J}_{d-1}$ it came from; this estimate is denoted by \hat{w} . Note that the left-to-right ordering of the nodes in \mathcal{S} and the original nodes in \mathcal{J}_{d-1} which

Algorithm 1 FINDPATHS in k -ary trees, TED deletion model

Input: a trace Y sampled from the TED deletion channel.

- 1: Initialize $\mathcal{S} = \emptyset$.
 - 2: **for** v a leaf in Y **do**
 - 3: **if** v has depth d in Y **then**
 - 4: Add the parent of v to \mathcal{S} .
 - 5: **end if**
 - 6: **end for**
 - 7: Initialize $\widehat{\mathcal{S}} = \emptyset$.
 - 8: **for** $v \in \mathcal{S}$ **do**
 - 9: **for** $\ell = 0$ to $d - 2$ **do**
 - 10: Compute \widehat{a}_ℓ based on node-to-leaf anchor paths and combining plug-in estimators (see equation (3) for the final formula, the text for further details, and Figure 3 for an illustration).
 - 11: **end for**
 - 12: Set $\widehat{w} := \widehat{a}_{d-2}\widehat{a}_{d-3} \cdots \widehat{a}_0$ (written in base- k expansion).
 - 13: Add \widehat{w} to $\widehat{\mathcal{S}}$.
 - 14: **end for**
 - 15: **Output:** $\widehat{\mathcal{S}}$.
-

they come from are the same, so the algorithm needs only to output the set $\widehat{\mathcal{S}} := \{\widehat{w} : v \in \mathcal{S}\}$ (since the mapping between \mathcal{S} and $\widehat{\mathcal{S}}$ follows the left-to-right ordering).²

Given $v \in \mathcal{S}$, to compute the estimate \widehat{w} , we first observe that any node $w \in \mathcal{J}_{d-1}$ can be written in its base- k expansion,

$$w = a_{d-2}a_{d-3} \cdots a_0,$$

where $a_\ell \in \{0, 1, \dots, k - 1\}$ for $\ell \in \{0, 1, \dots, d - 2\}$. Thus in order to compute an estimate \widehat{w} , it suffices to compute an estimate \widehat{a}_ℓ of a_ℓ for every $\ell \in \{0, 1, \dots, d - 2\}$ and then set

$$\widehat{w} := \widehat{a}_{d-2}\widehat{a}_{d-3} \cdots \widehat{a}_0.$$

The following is an equivalent and more pictorial way of thinking about this observation. Let u_0, u_1, \dots, u_{d-1} denote the nodes in X on the path from the root to $w \in \mathcal{J}_{d-1}$, with u_t having depth t (in particular, u_0 is the root and $u_{d-1} = w$). Then, for $\ell \in \{0, 1, \dots, d - 2\}$, the quantity a_ℓ is the position (from the left, with indexing starting at 0) of $u_{d-1-\ell}$ among the k children of $u_{d-2-\ell}$. Now suppose that $v \in \mathcal{S}$ came from node $w \in \mathcal{J}_{d-1}$, and let $u'_0, u'_1, \dots, u'_{d-1}$ denote the nodes in Y on the path from the root to v , with u'_t having depth t (in particular, u'_0 is the root and $u'_{d-1} = v$). Thus estimating a_0, a_1, \dots, a_{d-2} corresponds to estimating, for each node u'_t , where its pre-image in X ranks in the left-to-right ordering of itself and its $k - 1$ siblings.

We now explain how to compute the estimate \widehat{a}_0 given $v \in \mathcal{S}$; computing \widehat{a}_ℓ for general ℓ is similar but involves slightly more notation, so we defer this for now. Let z_0, z_1, \dots, z_m denote v and its siblings in Y , ordered from left to right, and let $\mathcal{Z} := \{z_0, z_1, \dots, z_m\}$. Let $z_0^*, z_1^*, \dots, z_{k'}^*$ denote the nodes among \mathcal{Z} that have a child in Y , ordered from left to right, and let $\mathcal{Z}^* := \{z_0^*, z_1^*, \dots, z_{k'}^*\}$. Note that $v \in \mathcal{Z}^*$ by definition; define k^* to be the index such that $v = z_{k^*}^*$. Also, by construction, the pre-images of all nodes in \mathcal{Z}^* were siblings in X , so we must have that $k' \leq k - 1$. Note that there are two ways that a node can be in $\mathcal{Z} \setminus \mathcal{Z}^*$:

²Regarding notation: note that $\widehat{\mathcal{S}}$ is *not* an estimate of \mathcal{S} , but rather an estimate of the pre-image of \mathcal{S} before X is passed through the deletion channel to obtain Y . We hope that the reader accepts this abuse of notational convention.

- A sibling w' of w in X is *not* deleted in Y , but *all* of the children of w' are deleted in Y . Then the image of w' in Y is in $\mathcal{Z} \setminus \mathcal{Z}^*$. Note that this is a highly unlikely event, since k is large.
- A sibling w' of w in X is deleted in Y , but not all of the children of w' are deleted in Y . Then the images of the nondeleted children of w' in Y are in $\mathcal{Z} \setminus \mathcal{Z}^*$. Note that if such a vertex w' is deleted, then in expectation there will be $(1 - q)k$ nondeleted children.

Since the first bullet point above is highly unlikely and the second bullet point describes the typical behavior of a trace, this motivates the following estimation procedure. For $i \in [k']$, let α_i denote the number of nodes in \mathcal{Z} that are between z_{i-1}^* and z_i^* ; furthermore, let α_0 denote the number of nodes in \mathcal{Z} that are before z_0^* . Now for every $i \in \{0, 1, \dots, k'\}$ let $\hat{\alpha}_i$ denote the unique integer satisfying

$$\hat{\alpha}_i - 1/2 \leq \frac{\alpha_i}{(1 - q)k} < \hat{\alpha}_i + 1/2.$$

Finally, we set

$$(1) \quad \hat{a}_0 := k^* + \sum_{i=0}^{k^*} \hat{\alpha}_i.$$

(If this results in an estimate that is greater than $k - 1$, then instead set $\hat{a}_0 := k - 1$.)

Now we turn to estimating \hat{a}_ℓ for general $\ell \in \{0, 1, \dots, d - 2\}$, given $v \in S$. Recall that $u'_0, u'_1, \dots, u'_{d-1}$ denote the nodes in Y on the path from the root to v , with u'_t having depth t . Let z_0, z_1, \dots, z_m denote $u'_{d-1-\ell}$ and its siblings in Y , ordered from left to right, and let $\mathcal{Z} := \{z_0, z_1, \dots, z_m\}$ (we reuse notation from above). Let $z_0^*, z_1^*, \dots, z_{k'}^*$ denote the nodes among \mathcal{Z} that have height $\ell + 1$ in Y , ordered from left to right, and let $\mathcal{Z}^* := \{z_0^*, z_1^*, \dots, z_{k'}^*\}$. Note that $u'_{d-1-\ell} \in \mathcal{Z}^*$ by definition; define k^* to be the index such that $u'_{d-1-\ell} = z_{k^*}^*$. Also, by construction, the pre-images of all nodes in \mathcal{Z}^* were siblings in X , so we must have that $k' \leq k - 1$. For $i \in [k']$, let α_i denote the number of nodes in Y that are either (a) in \mathcal{Z} between z_{i-1}^* and z_i^* , or (b) are descendants in Y of such a node; see Figure 3 for an illustration. Furthermore, let α_0 denote the number of nodes in Y that are either (a) in \mathcal{Z} before z_0^* , or (b) are descendants in Y of such a node. Now for every $i \in \{0, 1, \dots, k'\}$ let $\hat{\alpha}_i$ denote the unique integer satisfying

$$(2) \quad \hat{\alpha}_i - 1/2 \leq \frac{\alpha_i}{(1 - q) \sum_{h=1}^{\ell+1} k^h} < \hat{\alpha}_i + 1/2.$$

Finally, we again set

$$(3) \quad \hat{a}_\ell := k^* + \sum_{i=0}^{k^*} \hat{\alpha}_i.$$

The estimate in equation (3) is thus a generalization of the special case of \hat{a}_0 in equation (1). (If this results in an estimate that is greater than $k - 1$, then instead set $\hat{a}_\ell := k - 1$.)

This fully completes the description of the FINDPATHS algorithm. In the following lemma we analyze the performance of the FINDPATHS algorithm and show that its output is correct with high probability.

LEMMA 3.2. *There exist constants c and c' , that depend only on q , such that the following holds. Let $k \geq c \log^2(n)$, let X be a k -ary tree with arbitrary binary labels, and let Y be a trace sampled from the TED deletion channel. The FINDPATHS algorithm is fully successful—that is, for all nodes $v \in S$, the estimate \hat{w} is correct—with probability at least $1 - \exp(-c'\sqrt{k})$.*

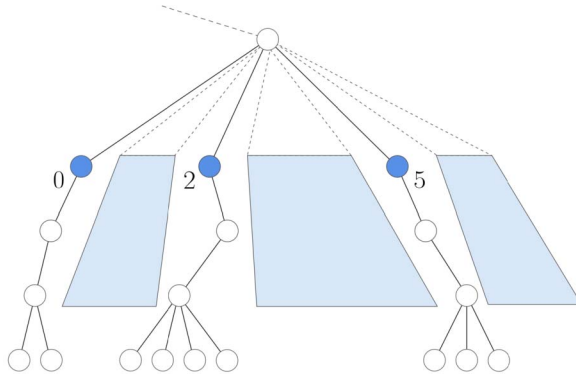


FIG. 3. Path estimation in k -ary trees. We estimate, level by level, the pre-image of each root-to-leaf path in a trace. At each level, we estimate the number of nodes deleted at that level using the number of nodes in the (light blue) trapezoids in the figure (this uses that k is large enough to apply concentration bounds). This then allows us to determine the positions of the surviving (dark blue) nodes, which have paths to a leaf (e.g., positions 0, 2, 5 above).

PROOF. Throughout this proof we denote the complement of an event \mathcal{E} by \mathcal{E}^c . Set $b := 10\sqrt{k}/\log(1/q) = \Omega(\log n)$ and let \mathcal{B}_b denote the event that Y is b -balanced. By Claim 1 we have that

$$(4) \quad \mathbb{P}(\mathcal{B}_b^c) \leq \exp(-C'\sqrt{k})$$

for some constant C' .

Let u be an internal node of X and let h_u be the height of u . Since u is an internal node, we have that $1 \leq h_u \leq d$. The number of descendants of u in X is $\sum_{\ell=1}^{h_u} k^\ell$. Let R_u denote the number of descendants of u in X that survive in Y .

Now fix $m \leq b$ and let u_1, \dots, u_m denote m consecutive siblings in X with height $h \in [d]$. Let $\mathcal{E}_{u_1, \dots, u_m}$ denote the event that

$$\left| \sum_{i=1}^m R_{u_i} - m(1-q) \sum_{\ell=1}^h k^\ell \right| \leq \frac{1}{3} (1-q) \sum_{\ell=1}^h k^\ell.$$

By a standard Chernoff bound we have that there exist constants $c'', c''' > 0$ such that

$$(5) \quad \mathbb{P}(\mathcal{E}_{u_1, \dots, u_m}^c) \leq \exp\left(-c''(1-q) \sum_{\ell=1}^h k^\ell / m\right) \leq \exp(-c''(1-q)k/m) \leq \exp(-c'''\sqrt{k}),$$

where in the last inequality we used that $m \leq b = 10\sqrt{k}/\log(1/q)$.

Finally, define the event

$$\mathcal{E} := \mathcal{B}_b \cap \bigcap_{m=1}^b \bigcap_{u_1, \dots, u_m} \mathcal{E}_{u_1, \dots, u_m},$$

where the intersection is over all possible m consecutive siblings u_1, \dots, u_m in X . Putting together equation (4), equation (5), and a union bound, we have that

$$\mathbb{P}(\mathcal{E}^c) \leq \exp(-c'\sqrt{k})$$

for some constant c' . On the other hand, on the event \mathcal{E} , the estimates $\hat{\alpha}_i$ in equation (2) are correct for all $v \in \mathcal{S}$, all $\ell \in \{0, 1, \dots, d-2\}$, and all $i \in \{0, 1, \dots, k'\}$. This implies that for every $v \in \mathcal{S}$ and every $\ell \in \{0, 1, \dots, d-2\}$, the estimate \hat{a}_ℓ in equation (3) is correct. Therefore for every $v \in \mathcal{S}$ the estimate \hat{w} is also correct. \square

Algorithm 2 Reconstructing k -ary trees, $k \geq c \log^2(n)$, TED deletion model

Set $T = \exp(c' \log_k(n)) \cdot T(k, 1/n^2)$ (for a large enough constant c').

Input: traces Y_1, \dots, Y_T sampled independently from the TED deletion channel.

```

1: for  $j \in \mathcal{J}_{d-1}$  do
2:   Initialize  $\mathcal{A}_j = \emptyset$ .
3: end for
4: for  $t = 1$  to  $T$  do
5:   Run Algorithm 1 with input  $Y_t$ ; let  $\widehat{\mathcal{S}}_t$  denote the output.
6:   for  $j \in \widehat{\mathcal{S}}_t$  do
7:     Add  $Y_t$  to  $\mathcal{A}_j$ .
8:   end for
9: end for
10: for  $j \in \mathcal{J}_{d-1}$  do estimate the labels of  $G_X(j)$  as follows:
11:   if  $\mathcal{A}_j = \emptyset$  then ▷ This happens with vanishing probability.
12:     Terminate the algorithm and produce no output.
13:   end if
   To estimate the labels of  $P_X(j)$ :
14:   Choose an arbitrary trace  $Y \in \mathcal{A}_j$ ;
15:   Let  $v$  denote the node in  $Y$  which caused  $Y$  to be included in  $\mathcal{A}_j$ ;
16:   Estimate labels of  $P_X(j)$  by copying bits from the path in  $Y$  that goes from the root
   to  $v$ .
   To estimate the labels of the children of  $j$ :
17:   Initialize  $\mathcal{T} = \emptyset$ .
18:   for  $Y \in \mathcal{A}_j$  do
19:     Let  $v$  denote the node in  $Y$  which caused  $Y$  to be included in  $\mathcal{A}_j$ ;
20:     Form a string  $Z$  by reading, from left to right, the bits of the children of  $v$  in  $Y$ ;
21:     Add  $Z$  to  $\mathcal{T}$ .
22:   end for
23:   Use a string trace reconstruction algorithm to estimate the labels of the children of  $j$ 
   from  $\mathcal{T}$ .
24: end for
25: Output: Take a union, over all  $j \in \mathcal{J}_{d-1}$ , of the estimated labels of  $G_X(j)$ , to estimate
   the labels of  $X$  (as in Remark 2).
```

The reconstruction algorithm: Estimating the labels of $G_X(j)$ for each $j \in \mathcal{J}_{d-1}$. Now that we have described and analyzed the FINDPATHS algorithm (Algorithm 1), we turn our attention to the full reconstruction algorithm (see Algorithm 2).

Lines 1–9 of Algorithm 2 describe the first step of the reconstruction algorithm, where we process all the traces and group them into different sets. Formally, we define a set \mathcal{A}_j for every $j \in \mathcal{J}_{d-1}$, which we initialize with $\mathcal{A}_j = \emptyset$. Then for every trace Y_t in our input, we run Algorithm 1 with input Y_t , and we let $\widehat{\mathcal{S}}_t$ denote the output. We then add Y_t to \mathcal{A}_j for every $j \in \widehat{\mathcal{S}}_t$.

We now turn to the main step of the reconstruction algorithm, which is described in lines 10–25 of Algorithm 2. For every $j \in \mathcal{J}_{d-1}$, we use the traces in \mathcal{A}_j to estimate the labels of $G_X(j)$, and finally we take a union of these estimates to estimate the labels of X . The estimation of the labels of $G_X(j)$ is done in two parts: (1) the estimation of the labels of $P_X(j)$, and (2) the estimation of the labels of the children of j ; see Figure 4 for an illustration.

To estimate the labels of $P_X(j)$, we take an arbitrary trace $Y \in \mathcal{A}_j$; if $\mathcal{A}_j = \emptyset$, then the algorithm terminates without output. Let v denote the node in Y which caused Y to be included

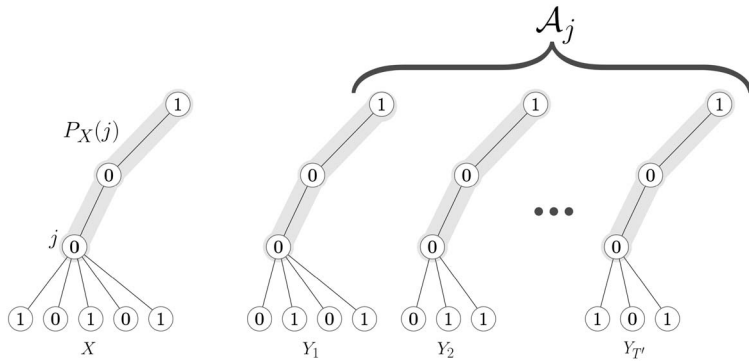


FIG. 4. An example set of traces \mathcal{A}_j appearing in Algorithm 2. The two salient points are that, with high probability: (1) the bits on the highlighted paths are all the same as the original bits in $P_X(j)$, and (2) the restriction to leaves corresponds to string trace reconstruction.

in \mathcal{A}_j . Assuming that Y was included in \mathcal{A}_j for the correct reason, that is, the pre-image of v is indeed j , then the labels of $P_X(j)$ are identical to the bits on the path in Y that goes from the root to v ; see Figure 4 for an illustration. Therefore we estimate the labels of $P_X(j)$ by copying the bits from the path in Y that goes from the root to v .

Finally, we estimate the labels of the children of j ; it turns out that this reduces to string trace reconstruction. Given a trace $Y \in \mathcal{A}_j$, let v denote the node in Y which caused Y to be included in \mathcal{A}_j . Assuming that Y was included in \mathcal{A}_j for the correct reason, that is, the pre-image of v is indeed j , then the children of v in Y are a random subset of the children of j in X ; see Figure 4 for an illustration. Thus if we restrict our attention to the bits on the children of j in X , the children of v in the trace Y are as if the original bits were passed through the string deletion channel; see Figure 4 again for an illustration. This motivates collecting a string trace from each $Y \in \mathcal{A}_j$, by looking at the children of the appropriate vertex v ; we let \mathcal{T} denote this collection of string traces. Finally, we use a string trace reconstruction algorithm to reconstruct the bits on the children of j in X from \mathcal{T} .

Now that we have fully described the reconstruction algorithm, we are ready to prove that it correctly reconstructs the labels of X with high probability. The following lemma is an important step towards this.

LEMMA 3.3. *There exist finite positive constants c' and c'' such that the following holds. Let $T = \exp(c' \log_k(n)) \cdot T(k, 1/n^2)$ and let Y_1, \dots, Y_T be i.i.d. traces from the TED deletion channel. With probability at least $1 - \exp(-c'' \sqrt{k})$ the following hold:*

1. The FINDPATHS algorithm is fully correct for all traces Y_1, \dots, Y_T .
2. For every $j \in \mathcal{J}_{d-1}$ we have that $|\mathcal{A}_j| \geq T(k, 1/n^2)$.

PROOF. The first claim follows from Lemma 3.2 and a union bound, using the fact that $T = \exp(O(\log_k(n) + k^{1/3}))$. For each $j \in \mathcal{J}_{d-1}$, the path $P_X(j)$ consists of $d - 1$ nonroot nodes and hence it survives in a trace with probability $(1 - q)^{d-1}$. Since $T \geq 2(1 - q)^{-(d-1)}T(k, 1/n^2)$, the second claim follows from a standard Chernoff bound. \square

Finishing the proof of Theorem 1.3.

PROOF OF THEOREM 1.3. The reconstruction algorithm is described in Algorithm 2, with a subroutine described in Algorithm 1. Let \mathcal{E} be the event that (1) the FINDPATHS algorithm is fully correct for all traces Y_1, \dots, Y_T , and (2) for every $j \in \mathcal{J}_{d-1}$ we have that $|\mathcal{A}_j| \geq T(k, 1/n^2)$. By Lemma 3.3 we have that $\mathbb{P}(\mathcal{E}) \geq 1 - \exp(-c'' \sqrt{k})$.

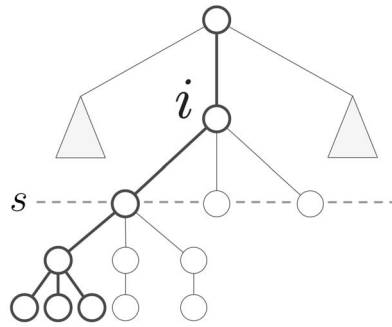


FIG. 5. An example of a trace that is s -stable for i (where here $s = 3$).

Conditioned on the event \mathcal{E} , the reconstruction algorithm correctly reconstructs the labels of $P_X(j)$ for every $j \in \mathcal{J}_{d-1}$ (see lines 14–17 of Algorithm 2); in other words, the reconstruction algorithm correctly reconstructs the labels of all internal nodes of X .

We next turn to the leaves of X . Conditioned on the event \mathcal{E} we have that $|\mathcal{A}_j| \geq T(k, 1/n^2)$, so using string trace reconstruction we can correctly reconstruct the labels of all children of j with probability at least $1 - 1/n^2$. Since there are at most n nodes in \mathcal{J}_{d-1} , a union bound shows that, conditioned on the event \mathcal{E} , we can correctly reconstruct the labels of all leaves of X with probability at least $1 - 1/n$.

Overall, the error probability in reconstructing the labels of X is at most $\exp(-c''\sqrt{k}) + 2/n$. \square

3.2. *Proof of Theorem 1.4 concerning arbitrary degree trees.* Recall the definition for \mathcal{I} defined in the first paragraph of Section 2.1: $\mathcal{I} := \bigcup_{t=1}^{d-1} \mathcal{I}_t$, where $\mathcal{I}_1 := \mathcal{J}_1 = \{0, 1, \dots, k-1\}$ and for $t \geq 2$, $\mathcal{I}_t := \{i \in \mathcal{J}_t \mid i \bmod k \neq 0\}$, and \mathcal{J}_t is the set of the nodes at depth t . We use traces that have a strong underlying structure, which we call s -stable; see Figure 5 for an illustration.

DEFINITION 3.4. A trace Y is s -stable for $i \in \mathcal{I}$ if $G_Y(i) \neq \perp$, and for every internal node v in $G_Y(i)$ with height $h \leq s$ in Y , each of the k children of v has height exactly $h - 1$ in Y .

Algorithm 3 below states, in pseudocode, our reconstruction algorithm for proving Theorem 1.4. At a high level, we will recover the labels for $G_X(i)$ separately for each $i \in \mathcal{I}$, which is sufficient because these subtrees cover all of the nonroot nodes in X .

The challenge is that, in the TED deletion model, $G_X(i)$ may shift to an incorrect position, even when $G_Y(i) \neq \perp$. This happens, for example, when the parent of i has children deleted in such a way that i moves to the left or right, but i still has $k - 1$ siblings (some of which are new); see Figures 6 and 7 for an illustration. The intuition for overcoming this issue is as follows. Let u be a node in $G_X(i)$ with child u' that is not a leaf (so u and u' both originally have k children). If u and all of its k children survive in a trace, then we will be in good shape. However, consider the situation when u survives and u' is deleted. In the TED model, we expect $(1 - q)k$ children of u' to move up to become children of u . Since this occurs for every deleted child of u , we expect u to now have many more than k children.

The bad case is when u has exactly k children in a trace after some of its original children are deleted; see Figure 7 for an illustration. This only happens when subtrees rooted at children of u are completely deleted. If such a subtree is large (i.e., u is higher up in the tree), then this is extremely unlikely. To deal with the nodes u closer to the leaves, we use the s -stable property to force the relevant subtrees to survive.

Algorithm 3 Reconstructing k -ary trees, arbitrary k , TED deletion model

Set $s = \lceil \log_k \log_{1/q}(3dk) \rceil$ and $T = C \log(n) \cdot (1 - q)^{-(dk+s^2k)}$ (for a large enough C).

Input: traces Y_1, \dots, Y_T sampled independently from the TED deletion channel.

- 1: Set $\mathcal{A} = \{Y_1, \dots, Y_T\}$.
 - 2: **for** $i \in \mathcal{I}$ **do**
 - 3: Initialize $\mathcal{A}_i = \emptyset$.
 - 4: **for** $t = 1$ to T **do**
 - 5: **if** Y_t is s -stable for i **then** add Y_t to \mathcal{A}_i .
 - 6: **end if**
 - 7: **end for**
 - 8: **for** node b in $G_X(i)$ **do**
 - 9: Let the learned label of b be the majority vote over all $Y \in \mathcal{A}_i$ of the labels on node b in $G_Y(i)$, as in Lemma 3.6.
 - 10: **end for**
 - 11: **end for**
 - 12: **Output:** Union the learned labels of $G_X(i)$ over all $i \in \mathcal{I}$ to reconstruct labels of X , as described in Remark 2.
-

An obvious way for Y to be s -stable is for it to contain $G_X(i)$ and enough relevant descendants of nodes in $G_X(i)$. Let $G_X^+(i)$ be the union of $G_X(i)$ and the k children of every internal node in $G_X(i)$; see Figure 6 for an illustration. Then Y will be s -stable if it contains $G_X^+(i)$ and at least one path to a leaf (in X) from every node in $G_X^+(i)$ with height at most s . In Lemma 3.5, we even argue that this happens with high enough probability to achieve the bound in the theorem.

Unfortunately, we cannot directly check whether Y contains the exact nodes in $G_X^+(i)$. We can check if Y is s -stable for i by examining the nodes of $G_Y(i)$ and their descendants in Y . But if Y is s -stable, then it is still not necessarily the case that $G_Y(i) = G_X(i)$, since the nodes in $G_X(i)$ may have shifted in Y or been deleted.

To get around this complication, we rely on the s -stable property of a trace. We argue in Lemma 3.6 that if s is large enough and a trace Y is s -stable for i , then with probability at least $2/3$, we have $G_Y(i) = G_X(i)$. We take a majority vote of $G_Y(i)$ over $O(\log n)$ traces Y to recover $G_X(i)$ with high probability. Since the subtrees $G_X(i)$ for $i \in \mathcal{I}$ cover X , we will be done.

Analyzing and using stable traces. In what follows, we fix $s = \lceil \log_k \log_{1/q}(3dk) \rceil$. We first show that a trace is s -stable with good enough probability.

LEMMA 3.5. *For $i \in \mathcal{I}$, a trace is s -stable for i with probability at least $(1 - q)^{dk+s^2k}$.*

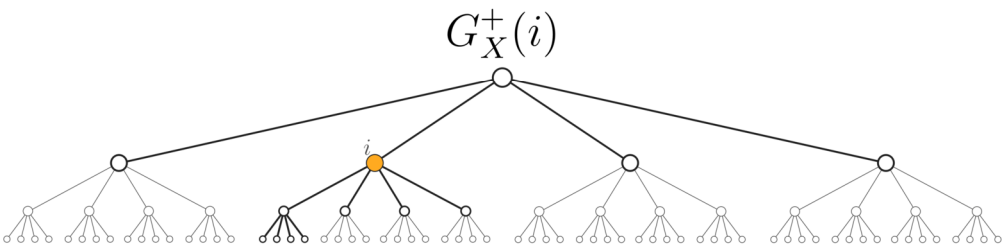


FIG. 6. Example of $G_X^+(i)$, where the node i is orange, and the full subtree is bold.

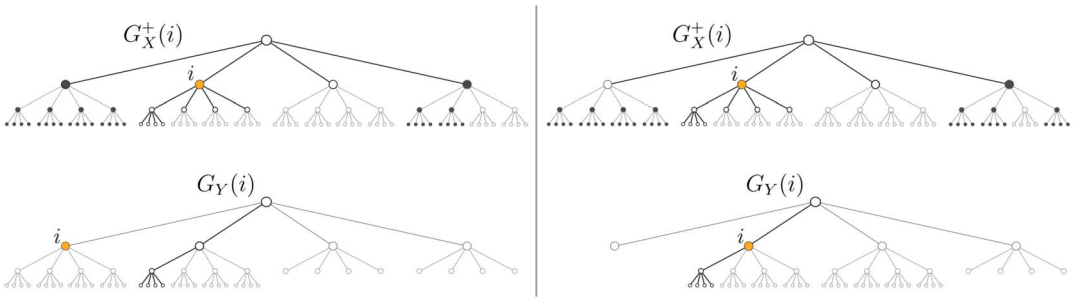


FIG. 7. Two traces, one “bad” and one “good”. In the top two trees, gray nodes in X are deleted to produce the corresponding traces below. The trace on the left has the subtree rooted at i in the incorrect location (it moved over to the left). The trace on the right has the subtree in the correct location.

PROOF. Being s -stable has two conditions. First, we need $G_Y(i) \neq \perp$. Let $G_X^+(i)$ be the union of $G_X(i)$ and the k children of every internal node in $G_X(i)$, where $|G_X^+(i)| = dk + 1$. We will prove that if Y contains $G_X^+(i)$, then $G_Y(i) \neq \perp$, because in fact, $G_Y(i) = G_X(i)$. Since the root is never deleted, all nodes in $G_X^+(i)$ survive in a trace with probability $(1 - q)^{dk}$, and so $G_Y(i) = G_X(i)$ with at least this probability.

Assume that Y contains $G_X^+(i)$. Let $G_X(i) = u_0, \dots, u_{d+k-1}$, and consider building $G_Y(i) = v_0, \dots, v_{d+k-1}$ using π_i . We argue recursively: For $t \in [d - 1]$, we assume that $v_{t'} = u_{t'}$ for all $t' < t$, and we prove that $v_t = u_t$ as well. The base case $t' = 0$ holds because the root $v_0 = u_0$ is never deleted. Then, since Y contains $G_X^+(i)$, we know that $v_{t'}$ has exactly k children in Y , which are the children of $u_{t'}$ in X . Moreover, the left-to-right order of these k children is preserved in the deletion model. Therefore, the child of $v_{t'}$ in position $\pi_i(t')$ must indeed be $u_{t'+1}$ for all $t' < t$. This establishes $v_t = u_t$ for all $t \in \{0, 1, \dots, d - 1\}$. For the leaves of $G_X(i)$, when $v_{d-1} = u_{d-1}$, and v_{d-1} has k children in Y , then we must also have $v_d, \dots, v_{d+k-1} = u_d, \dots, u_{d+k-1}$.

For the second condition of s -stability, consider an internal node u_t in $G_X(i)$ with height $h = d - t$ satisfying $1 \leq h \leq s$. Let u'_0, \dots, u'_{k-1} be the children of u_t in X . Because u'_j has height $h - 1$ in X , there is some path with h nodes from u'_j to a leaf in X . Consider one such path for each $j = 0, \dots, k - 1$ such that $j \neq \pi_i(t)$. Since there are $k - 1$ choices for j , let P_t be the union of these $k - 1$ paths, where $|P_t| = h(k - 1) \leq s(k - 1)$. The survival of P_t guarantees that u'_j has the correct height for Y to be s -stable. Since $|\bigcup_{t=d-s}^{d-1} P_t| \leq s^2(k - 1)$, and each node survives independently with probability $(1 - q)$, we have that P_{d-s}, \dots, P_{d-1} survive with probability at least $(1 - q)^{s^2(k-1)}$.

Combining these two conditions, Y is s -stable with probability at least $(1 - q)^{dk+s^2k}$. \square

We now formalize the intuition that if all nodes in $G_Y(i)$ have k children, and the parents are high enough in the tree, then the children are probably correct. The reason is that subtrees rooted at their children are unlikely to be completely deleted. This is the only bad case, since otherwise, we expect deleted nodes to cause their parents to have many more than k children. Finally, since the trace is s -stable, the nodes near the leaves will be correct as well.

LEMMA 3.6. For $i \in \mathcal{I}$, if Y is a random s -stable trace for i , then $G_Y(i) = G_X(i)$ with probability at least $2/3$.

PROOF. Since Y is s -stable, $G_Y(i) \neq \perp$. Let $G_Y(i) = v_0, \dots, v_{d+k-1}$ and $G_X(i) = u_0, \dots, u_{d+k-1}$, where v_t and u_t have depth $t \in \{0, 1, \dots, d - 1\}$, and v_{d-1} and u_{d-1} have children v_d, \dots, v_{d+k-1} and u_d, \dots, u_{d+k-1} , respectively. Our strategy is to define an event

\mathcal{E} that happens with probability at least $2/3$ and implies that $v_t = u_t$ for $t \leq d + k - 1$. Consider $t \in [d]$, and let u'_0, \dots, u'_{k-1} be the children of u_{t-1} in X . Define \mathcal{E}_t to be the event that, for every $j \in \{0, 1, \dots, k - 1\}$, at least one node in the subtree rooted at u'_j survives in Y . Then, define $\mathcal{E}_{\leq m} = \bigcap_{t=1}^m \mathcal{E}_t$ and set $\mathcal{E} = \mathcal{E}_{\leq d}$.

We first argue that when $\mathcal{E}_{\leq m}$ holds, then $v_t = u_t$ for all $t \leq m$. Because the root has not been deleted, we have $v_0 = u_0$. Then, for $t \in [m]$, we assume that $v_{t'} = u_{t'}$ for $t' < t$, and we prove that $v_t = u_t$.

Because Y is s -stable, v_{t-1} has k children in Y . Denote them v'_0, \dots, v'_{k-1} . We need to show that u_t is in position $\pi_i(t - 1)$ among them, so that $v_t = v'_{\pi_i(t-1)} = u_t$. Since \mathcal{E}_t holds, there is some surviving node in Y from the subtree rooted at each original child of u_{t-1} in X . Moreover, since $u_{t-1} = v_{t-1}$, this accounts for at least k children of v_{t-1} in Y . Because there are exactly k children of v_{t-1} , it must be the case that $v'_{\pi_i(t-1)}$ is originally from the subtree rooted at u_t in X . In particular, $v'_{\pi_i(t-1)} = u_t$ if and only if u_t survives in Y .

We claim that if u_t were deleted, then it would contradict Y being s -stable, since we would have $G_Y(i) = \perp$ instead. Indeed, the deletion of u_t would cause $v'_{\pi_i(t-1)}$ to have height less than $d - t$ in Y . This would imply that at some depth d' with $t < d' < d$, the node $v_{d'}$ in $G_Y(i)$ would be a leaf, leading to $G_Y(i) = \perp$. We conclude that u_t survives in Y , and so that $v_t = v'_{\pi_i(t-1)} = u_t$, as desired.

We have shown that \mathcal{E} guarantees that $v_t = u_t$ for all $t \leq d - 1$. In particular, $v_{d-1} = u_{d-1}$, and the k children of v_{d-1} in Y must be the children of u_{d-1} in X . This finishes the argument that \mathcal{E} implies that $v_t = u_t$ for all $t \leq d + k - 1$, that is, $G_Y(i) = G_X(i)$.

Now we prove that \mathcal{E} happens with probability at least $2/3$ in an s -stable trace. We prove this in two steps. First, we argue that $\mathcal{E}_{\leq d-s}$ occurs with probability at least $2/3$. Then, we show that $\mathcal{E}_{\leq d-s}$ implies \mathcal{E} . Consider the node u_{t-1} in $G_X(i)$ for $t \in [d - s]$, and let u'_0, \dots, u'_{k-1} be the k children of u_{t-1} in X . Since the height of u'_j is at least s , the subtree rooted at u'_j in X contains at least $\sum_{\ell=0}^s k^\ell \geq k^s$ nodes. The probability that all of these nodes are deleted is at most q^{k^s} . Because $s = \lceil \log_k \log_{1/q}(3dk) \rceil$, this is at most $1/(3dk)$. Taking a union bound over the k children implies that \mathcal{E}_t occurs with probability at least $1 - 1/(3d)$, and taking a union bound over $t \in [d - s]$ implies that $\mathcal{E}_{\leq d-s}$ holds with probability at least $2/3$.

The final step is to prove that \mathcal{E} happens with probability one, in an s -stable trace, assuming that $\mathcal{E}_{\leq d-s}$ holds. More precisely, we will show that $\mathcal{E}_{\leq d-s+\ell}$ implies $\mathcal{E}_{d-s+\ell+1}$ for $\ell = 0, 1, \dots, s - 1$. We have already argued that $\mathcal{E}_{\leq d-s+\ell}$ guarantees that $v_{d-s+\ell} = u_{d-s+\ell}$. We claim that the k children v'_0, \dots, v'_{k-1} of $v_{d-s+\ell}$ are the original children of $u_{d-s+\ell}$ in X (and this clearly implies $\mathcal{E}_{d-s+\ell+1}$). Since Y is s -stable, there is a path with $s - \ell + 1$ nodes from v'_j to a leaf in Y . If v'_j were not an original child of $u_{d-s+\ell}$, then all such paths would have at most $s - \ell$ nodes. This implies no children of $u_{d-s+\ell} = v_{d-s+\ell}$ have been deleted in Y , and their existence witnesses the survival of the subtrees needed for $\mathcal{E}_{d-s+\ell+1}$. Since this holds for $\ell = 0, 1, \dots, s$, we conclude that $\mathcal{E} = \mathcal{E}_{\leq d}$ follows from $\mathcal{E}_{\leq d-s}$ in an s -stable trace, and $\Pr[G_Y(i) = G_X(i)] \geq \Pr[\mathcal{E}] = \Pr[\mathcal{E}_{\leq d-s}] \geq 2/3$. \square

Completing the proof of Theorem 1.4.

PROOF OF THEOREM 1.4. Let \mathcal{A} be a set of $T = C \log(n)/(1 - q)^{dk+s^2k}$ traces with C a large enough constant. By Lemma 3.5, each trace in \mathcal{A} is s -stable for i with probability $(1 - q)^{dk+s^2k}$. Therefore, by setting C large enough and taking a union bound over $i \in \mathcal{I}$, we can ensure that with probability at least $1 - 1/n^2$, for every $i \in \mathcal{I}$ there is a subset $\mathcal{A}_i \subseteq \mathcal{A}$ of s -stable traces for i with $|\mathcal{A}_i| \geq C' \log n$, for a constant C' to be set later.

By Lemma 3.6, each trace $Y \in \mathcal{A}_i$ has the property that $G_Y(i) = G_X(i)$ with probability at least $2/3$. Let $f_i(Y) \in \{0, 1\}^{d+k-1}$ be the labels of $G_Y(i)$ in Y . In expectation over $Y \in \mathcal{A}_i$, we

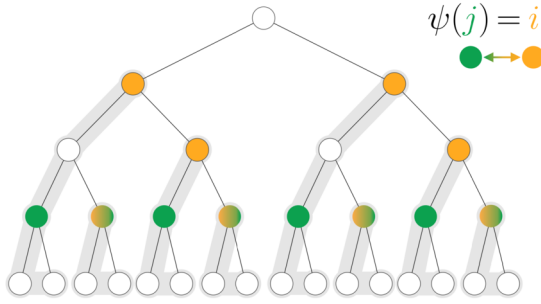


FIG. 8. The partition of the tree X into the subtrees $\{H_X(i) : i \in \mathcal{I}\}$ is illustrated in gray. The bijection $\psi : \mathcal{J}_{d-1} \rightarrow \mathcal{I}$ maps, for each gray subtree, a single green node to a single orange node in that subtree. Nodes that are both green and orange map to themselves.

have that at least a $2/3$ fraction of Y satisfy $f_i(Y) = f_i(X)$. Therefore, since $|\mathcal{A}_i| \geq C' \log n$ for a large enough constant C' , we have by a standard Chernoff bound that the majority value of $f_i(Y)$ over $Y \in \mathcal{A}_i$ is equal to $f_i(X)$, with probability at least $1 - 1/n^2$. For each $i \in \mathcal{I}$, our reconstruction algorithm uses this majority vote to deduce the labels for $G_X(i)$. Taking a union bound over $i \in \mathcal{I}$, where $|\mathcal{I}| \leq n$, we correctly label all nodes with probability at least $1 - 1/n$.

It remains to show that $T = \exp(O(dk))$, where $d = O(\log_k n)$. Recall that we have set $s = \lceil \log_k \log_{1/q}(3dk) \rceil$. If $k \geq d$, then $s \leq c \log \log k / \log k$ for a constant c , since q is a constant, and so $s = O(1)$. If $k \leq d$, then $s \leq c \log \log d$, and in particular, $s^2 < c''d$ for some constant c'' . Therefore, for any k , we have $T \leq C \log n \cdot \exp(c'dk)$ for some constant $c' > 1$ depending only on q , and since $\log \log n < dk$, we conclude that that $T = \exp(O(dk))$. \square

4. Reconstructing trees, left-propagation model. In this section we present our two algorithms for k -ary trees in the left-propagation deletion model.

4.1. *Proof of Theorem 1.5 concerning large degree trees.* Recall the definitions from the first paragraph in Section 2.1, in particular that of \mathcal{I} . Recall also that the sets $H_X(i)$ for $i \in \mathcal{I}$ partition the nonroot nodes of X , and each $H_X(i)$ contains exactly one node from \mathcal{J}_{d-1} ; see Figure 8 for an illustration. Define the bijection $\psi : \mathcal{J}_{d-1} \rightarrow \mathcal{I}$ as $\psi(j) = i$ for the distinct i such that $j \in H_X(i)$; see Figure 8 again for an illustration. For each fixed $j \in \mathcal{J}_{d-1}$ and each trace Y , we will extract a bit-string $s_Y(j)$ and use it to reconstruct the labels for $H_X(\psi(j))$. We only define $s_Y(j)$ whenever $P_Y(j) \neq \perp$, but this suffices for our purposes (as discussed later). To define $s_Y(j)$, we need some notation. Let v_0, v_1, \dots, v_{d-1} be the nodes in $P_Y(j)$, where v_t has depth t in Y . Then, let $v_d, \dots, v_{d+k'}$ be the children of v_{d-1} in Y , where $k' \leq k - 1$. Let $i = \psi(j)$ and let t_i be the depth of i in X . Finally, define $s_Y(j)$ as a bit-string of length $d + k' - t_i + 1$ consisting of the labels in Y of the nodes $v_{t_i}, \dots, v_{d+k'}$; see Figure 9 for an illustration.

Algorithm 4 below states, in pseudocode, our reconstruction algorithm for proving Theorem 1.5. We note that this reconstruction algorithm is tailored to the left-propagation deletion model.

For complete k -ary trees with sufficiently large $k \geq c \log n$, a trace Y has $P_Y(j) \neq \perp$ for all nodes $j \in \mathcal{J}_{d-1}$ with high probability. When $P_Y(j) \neq \perp$ and $j \in H_X(i)$, we can extract a subset of $H_X(i)$ that behaves as if it went through the string deletion channel (i.e., as if it were a path on $|H_X(i)|$ nodes). Therefore, using traces with $P_Y(j) \neq \perp$ for all $j \in \mathcal{J}_{d-1}$, we reduce to string trace reconstruction (see Figure 9 for an illustration), and we reconstruct the labels for each $H_X(i)$ separately. This suffices because the subtrees $H_X(i)$ for $i \in \mathcal{I}$ partition the nonroot nodes of X (see Figure 8 for an illustration).

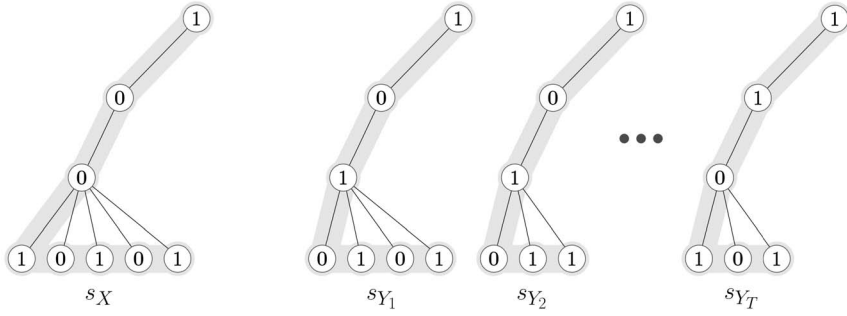


FIG. 9. Extracting the strings $s_Y(j)$ from the subtrees that contain the path $P_Y(j)$.

We first argue that $P_Y(j) \neq \perp$ with high probability when $k \geq c \log n$.

LEMMA 4.1. *Let $k \geq c \log n$. In the Left-Propagation model, a random trace Y has $P_Y(j) \neq \perp$ for every $j \in \mathcal{J}_{d-1}$ with probability at least $1 - \exp(-c'k)$.*

PROOF. The property that Y has $P_Y(j) \neq \perp$ for every $j \in \mathcal{J}_{d-1}$ is equivalent to Y containing a complete k -ary subtree of depth $d - 1$ with the same root as X . Consider any node $j \in \mathcal{J}_{d-1}$, and recall that the subtree $G_X(j)$ has $d + k - 1$ nonroot nodes. Each nonroot node in $G_X(j)$ survives independently in Y with probability $(1 - q)$. Let \mathcal{E}' be the event that at least d nodes from $G_X(j)$ survive in Y for every $j \in \mathcal{J}_{d-1}$. Because $k \geq c \log n$ and $d \leq \log_k n$ and $|\mathcal{J}_{d-1}| \leq n$, a standard Chernoff and union bound implies that \mathcal{E}' holds with probability $1 - \exp(-c'k)$ for a constant $c' > 0$ depending on q . When \mathcal{E}' holds, for all $j \in \mathcal{J}_{d-1}$, every node in $P_Y(j)$ has exactly k children in Y and $P_Y(j) \neq \perp$. \square

PROOF OF THEOREM 1.5. Let $T = T(d + k, 1/n^2)$ be the number of traces needed to learn $d + k$ bits with probability $1 - 1/n^2$ in the string model with deletion probability q . We will reconstruct X with probability $1 - O(1/n)$ using T traces from the left-propagation model.

By Lemma 4.1, a trace Y has $P_Y(j) \neq \perp$ for every $j \in \mathcal{J}_{d-1}$ with probability $1 - \exp(-c'k)$. By Theorem 1.1 we have that $T = \exp(O((d + k)^{1/3})) = \exp(O(k^{1/3}))$, where the second inequality is due to $d \leq k$. Thus by a union bound it follows that, with probability

Algorithm 4 Reconstructing k -ary trees, $k \geq c \log(n)$, left-propagation deletion model

Set $T = T(d + k, 1/n^2)$.

Input: traces Y_1, \dots, Y_T sampled independently from the TED deletion channel.

- 1: Set $\mathcal{A} = \{Y_1, \dots, Y_T\}$.
 - 2: **if** there exists a trace $Y \in \mathcal{A}$ and a node $j \in \mathcal{J}_{d-1}$ such that $P_Y(j) = \perp$ \triangleright This happens with vanishing probability.
 - 3: **then** Terminate the algorithm and produce no output.
 - 4: **else** \triangleright This happens with high probability.
 - 5: **for** $j \in \mathcal{J}_{d-1}$ **do**
 - 6: Reconstruct labels of $H_X(\psi(j))$ from $\{s_Y(j)\}_{Y \in \mathcal{A}}$, via string trace reconstruction (Theorem 1.1).
 - 7: **end for**
 - 8: **Output:** Union the learned labels of $H_X(i)$ over all $i \in \mathcal{I}$ to reconstruct labels of X , as described in Remark 2.
 - 9: **end if**
-

at least $1 - \exp(-c''k)$ for some constant $c'' > 0$, we have $P_Y(j) \neq \perp$ for all traces Y and nodes $j \in \mathcal{J}_{d-1}$. So from now on we assume that $P_Y(j) \neq \perp$ for every trace Y and every $j \in \mathcal{J}_{d-1}$.

Decompose X into subtrees $H_X(\psi(j))$ for $j \in \mathcal{J}_{d-1}$; see Figure 8 for an illustration. For each of the T traces Y , extract the bit-string $s_Y(j)$; see Figure 9 for an illustration. Consider these as T traces from the string deletion model on $|H_X(\psi(j))| < d + k$ bits. More precisely, let $s_X(j)$ be the labels in X for the nodes in $H_X(\psi(j))$. We claim that $s_Y(j)$ is a valid trace for the string deletion model with unknown string $s_X(j)$. In the left-propagation model, when $P_Y(j) \neq \perp$, the nodes considered in Y for $s_Y(j)$ form a subsequence of the corresponding nodes in X . Therefore, since each node is deleted with probability q , the bits in $s_Y(j)$ will be a trace of the string $s_X(j)$. Though we only consider traces with at least d bits remaining, the probability that at least one trace of T has less than d bits occurs with probability at most $\exp(-O(k))$. So by slightly increasing the factor C' in T , we can use Theorem 1.1 to see T traces suffice to reconstruct $s_X(j)$ with probability $1 - 1/n^2$. Moreover, $s_X(j)$ are the labels for $H_X(\psi(j))$. Taking a union bound over $|\mathcal{I}| \leq n$, we can reconstruct $H_X(i)$ for all $i \in \mathcal{I}$ with probability at least $1 - 1/n$. \square

4.2. Proof of Theorem 1.6 concerning arbitrary degree trees. Algorithm 5 below states, in pseudocode, our reconstruction algorithm for proving Theorem 1.6.

As in the proof of Theorem 1.5, we reconstruct X by reconstructing the subtrees $H_X(i)$ for $i \in \mathcal{I}$, which partition the nonroot nodes of X . Instead of reducing to string reconstruction, we use traces with $G_Y(i) \neq \perp$ to directly obtain labels for $H_X(i)$. We only need to take enough traces to balance out the fact that a trace with $G_Y(i) \neq \perp$ for $i \in \mathcal{I}$ occurs with probability $\exp(-O(d+k))$.

Recovering the labels for subtrees. We first show that if a trace satisfies $G_Y(i) \neq \perp$, then we can reconstruct the labels of $H_X(i)$; see Figure 10 for an illustration.

Algorithm 5 Reconstructing k -ary trees, arbitrary k , Left Propagation deletion model

Set $T = C(1 - q)^{-(d+c'k)} \log(n)$ (for large enough c' and C).

Input: traces Y_1, \dots, Y_T sampled independently from the TED deletion channel.

```

1: for every  $i \in \mathcal{I}$  do
2:   Initialize  $\mathcal{A}_i = \emptyset$ .
3:   for  $t = 1$  to  $T$  do
4:     if  $Y_t$  has  $G_{Y_t}(i) \neq \perp$  then add  $Y_t$  to  $\mathcal{A}_i$ .
5:     end if
6:   end for
7: end for
8: if there exists  $i \in \mathcal{I}$  such that  $\mathcal{A}_i = \emptyset$  ▷ This happens with vanishing probability.
9:   then Terminate the algorithm and produce no output.
10: else ▷ This happens with high probability.
11:   for every  $i \in \mathcal{I}$  do
12:     Choose an arbitrary  $Y \in \mathcal{A}_i$ .
13:     Reconstruct labels in  $H_X(i)$  bit by bit as those of  $H_Y(i)$ , using Remark 1.
14:   end for
15:   Output: Union the learned labels of  $H_X(i)$  over all  $i \in \mathcal{I}$  to reconstruct labels of  $X$ ,
    as described in Remark 2.
16: end if

```

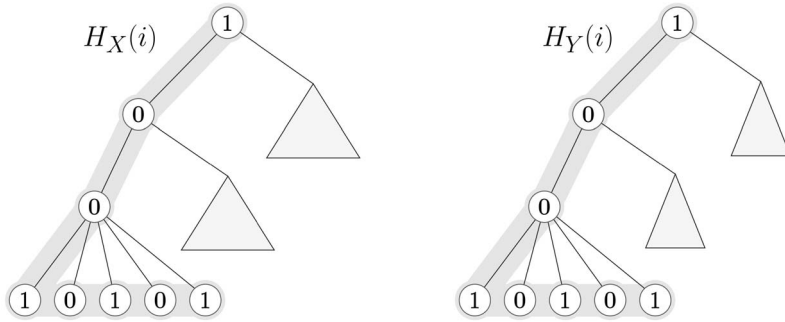


FIG. 10. Extracting correct labels for $s_X(j)$ from a single trace containing a caterpillar, $G_Y(i) \neq \perp$.

LEMMA 4.2. *In the left-propagation model, if $G_Y(i) \neq \perp$, then $H_Y(i) = H_X(i)$ and the labels for these subtrees are identical in Y and X .*

PROOF. Let $s_Y(i) = v_\ell, \dots, v_b$ be the labels on the left only path from i to the leaf plus its siblings in trace Y . These are the labels on $H_Y(i)$. Similarly, let $s_X(i) = u_{\ell'}, \dots, u_{b'}$ be the labels on the left-only path from i to the leaf plus its siblings in X . These are the labels on $H_X(i)$. Due to the behavior of the left-propagation model, the string $s_Y(i)$ is a trace of the string $s_X(i)$. When $G_Y(i) \neq \perp$, no nodes of $H_X(i)$ could have been deleted to obtain the trace $H_Y(i)$, otherwise the number of leaves in $G_Y(i)$ would be too small and $G_Y(i)$ would not be defined. As $s_Y(i)$ is a trace of $s_X(i)$ and $|H_Y(i)| = |H_X(i)|$, this implies that $H_Y(i) = H_X(i)$. \square

LEMMA 4.3. *Fix $i \in \mathcal{I}$ and let Y be a trace from the left-propagation deletion model. There exists an absolute constant $c' > 1$ such that with probability at least $(1 - q)^{d+c'k}$ we have that $G_Y(i) \neq \perp$.*

PROOF. There are $d + k$ nodes in $G_X(i)$ and they all survive in Y with probability $(1 - q)^{d+k}$. From now on, we assume this holds. Let $G_X(i) = u_0, \dots, u_{d+k-1}$, where u_t has depth $t \in \{0, 1, \dots, d\}$, and u_{d-1} has children u_d, \dots, u_{d+k-1} . Consider $t \in [d]$, and let u'_0, \dots, u'_{k-1} be the k children of u_{t-1} in X . Define \mathcal{E}_t to be the event that, for every $j \in \{0, 1, \dots, k - 1\}$, at least one node in the subtree rooted at u'_j survives in Y . We observe that, in the left-propagation model, if both $\bigcap_{t=1}^d \mathcal{E}_t$ holds and $G_X(i)$ survives, then we have $G_Y(i) \neq \perp$.

Because $G_X(i)$ surviving implies that the k children of u_{d-1} survive, we already know that \mathcal{E}_d holds. For $t \in [d - 1]$, each child u'_j of u_{t-1} has height $h = d - t + 1$ in X . In particular, the subtree rooted at u'_j in X contains at least k^{d-t+1} nodes. If $u'_j \in G_X(i)$, then we have assumed it survives, otherwise there are $k - 1$ other subtrees. Since the subtrees considered are independent, at least one node survives from each of them (for all $t \in [d - 1]$) with probability at least

$$\prod_{h=2}^d (1 - q^{k^h})^{k-1} = (1 - q)^{ck},$$

for some constant c . Putting everything together, $G_X(i)$ survives and $\bigcap_{t=1}^d \mathcal{E}_t$ holds, and therefore $G_Y(i) \neq \perp$, with probability at least $(1 - q)^{d+k} \cdot (1 - q)^{ck} = (1 - q)^{d+(c+1)k}$. \square

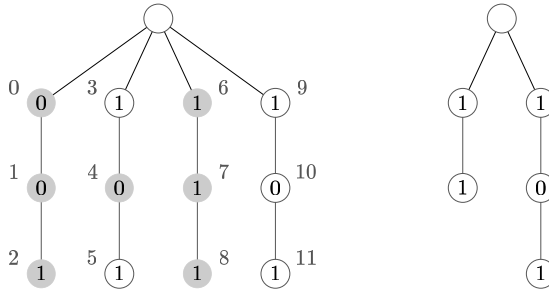


FIG. 11. DFS indexing and example trace (in both deletion models) for a (12, 3)-spider.

Completing the proof of Theorem 1.6.

PROOF OF THEOREM 1.6. By Lemma 4.3, the probability that none of T traces satisfy $G_Y(i) \neq \perp$ for some $i \in \mathcal{I}$ is at most

$$|\mathcal{I}|(1 - (1 - q)^{d+c'k})^T \leq n \exp(-T(1 - q)^{d+c'k}).$$

To ensure that at least one trace has $G_Y(i) \neq \perp$ for every $i \in \mathcal{I}$ with high probability, we take $T' = O(T \log n)$ traces, where T satisfies

$$T = (1 - q)^{-(d+c'k)} \leq (1 - q)^{-\ln(n)/\ln(k) - c'k} = n^{\ln(1/(1-q))(c'k/\ln(n)+1/\ln(k))}.$$

By Lemma 4.2, any trace with $G_Y(i) \neq \perp$ induces the correct labeling of $H_X(i)$ by using the labels for the nodes in $H_Y(i)$. In other words, with high probability, a set of $T' = O(T \log n)$ traces yields a correct labeling of $H_X(i)$ for all $i \in \mathcal{I}$. Since the subtrees $H_X(i)$ for $i \in \mathcal{I}$ form a partition of X , we can recover all labels in X . \square

5. Reconstructing spiders. In this section, we describe how to reconstruct spiders and prove Theorem 1.7 and Proposition 1.8. We start with preliminaries in Section 5.1. An outline of the proof for Theorem 1.7 is followed by the full proof in Section 5.2. The proof assumes a lemma requiring complex analysis that is deferred to Section 5.3. Proposition 1.8 is proven in Section 5.4. The remaining proofs of lemmas stated in this section are detailed in Section 5.5.

5.1. *Spider algorithm preliminaries.* When a labeled (n, d) -spider, X , goes through the deletion channel, we assume that its trace, Y , is an (n, d) -spider by inserting length d paths of 0 s after the remaining paths and nodes labeled 0 to the end of paths. After this, traces have n/d paths of length d (excluding the root).

We define a left-to-right ordered DFS index for (n, d) -spiders, illustrated in Figure 11. The labels increase along the length of the paths from the root and increase left to right among the paths. Specifically, if node v is in the i th path from the left and has depth j , then its label is $(i - 1)d + j - 1$. These labels will be used to define appropriate generating functions. As discussed in Remark 4, we need not consider the root as part of the generating function.

5.2. *Proof of Theorem 1.7 concerning (n, d) -spiders with small d .* In the regime where spiders have short paths ($d \leq \log_{1/q} n$), we use mean-based algorithms that generalize the methods of [14, 15, 34]. Using the DFS indexing of nodes, let X be an (n, d) -spider with labels $\{a_j\}_{j=0}^{n-1}$ and let Y be a trace of X , with the labels of Y denoted by $\{b_j\}_{j=0}^{n-1}$. Consider now the random generating function

$$\sum_{j=0}^{n-1} b_j w^j$$

for $w \in \mathbb{C}$. Due to the special structure of spiders, the expected value of this random generating function can be computed (see Lemma 5.1 below), and while it is more complicated than the corresponding formula for strings, it is still tractable. This is useful since by averaging samples we can approximate this expected value.

We then show that for every pair of labeled (n, d) -spiders, X^1 and X^2 , with different binary labels, we can carefully choose $w \in \mathbb{C}$ so that the corresponding values of the respective generating functions differ in expectation at some index $j = j(X^1, X^2)$. In choosing between candidate spiders X^1 and X^2 , the algorithm deems the better match of the pair to be the spider for which the expected value of the generating function at $w \in \mathbb{C}$ is closer to the mean of the traces at j . If any spider is a better match compared to every other spider, it is said to be the best match, and the algorithm outputs that spider.

For the quantitative estimates, the key technical challenge is to lower bound the modulus of the (expected) generating function on a carefully chosen arc of the unit disc in the complex plane. Our analysis, based on harmonic measure, is inspired by [8], as well as the recent work of [20].

When d is constant, the reconstruction problem on (n, d) -spiders can be reduced to string trace reconstruction (see Proposition 5.8). Hence, we will assume that d is greater than a specific constant ($d \geq 20$ suffices). We begin by computing the expected value of the generating function for an (n, d) -spider which has gone through a deletion channel with parameter q . We denote this expected generating function by $A(w)$, where $w \in \mathbb{C}$.

LEMMA 5.1. *Let $a = \{a_i\}_{i=0}^{n-1}$ be the labels of an (n, d) -spider with labels $a_i \in \mathbb{R}$ and let $b = \{b_j\}_{j=0}^{n-1}$ be the labels of its trace from the deletion channel with deletion probability q . Then*

$$A(w) := \mathbb{E} \left(\sum_{j=0}^{n-1} b_j w^j \right) = (1 - q) \sum_{\ell=0}^{n-1} a_\ell (q + (1 - q)w)^\ell \pmod{d} (q^d + (1 - q^d)w^d)^{\lfloor \frac{\ell}{d} \rfloor},$$

where the expectation is over the random labels b .

While $A(w)$ is written as only a function of w , it implicitly depends on the labels a of the original spider. The proof of Lemma 5.1 is in Section 5.5, as it follows from a standard manipulation of equations. We use this generating function to distinguish between two candidate (n, d) -spiders X^1 and X^2 , which have labels $a^1 = \{a_j^1\}_{j=0}^{n-1}$ and $a^2 = \{a_j^2\}_{j=0}^{n-1}$ which are different (i.e., there exists $j \in \{0, 1, \dots, n - 1\}$ such that $a_j^1 \neq a_j^2$). Let Y^1 and Y^2 denote random traces with labels $b^1 = \{b_j^1\}_{j=0}^{n-1}$ and $b^2 = \{b_j^2\}_{j=0}^{n-1}$ that arise from passing X^1 and X^2 through the deletion channel with deletion probability q .

Define $a := a^1 - a^2$ and let $A(w)$ be the expected generating function with input a . From Lemma 5.1 we have that

$$(6) \quad \sum_{j=0}^{n-1} (\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]) w^j = A(w).$$

Let $\ell^* := \arg \min_{\ell \geq 0} \{a_\ell \neq 0\}$ (note that $\ell^* \leq n - 1$ by construction) and define

$$\tilde{A}(w) := (1 - q) \sum_{\ell=\ell^*}^{n-1} a_\ell (q + (1 - q)w)^\ell \pmod{d} (q^d + (1 - q^d)w^d)^{\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor}.$$

Observe that $A(w) = (q^d + (1 - q^d)w^d)^{\lfloor \frac{\ell^*}{d} \rfloor} \cdot \tilde{A}(w)$; accordingly, we call $\tilde{A}(w)$ the factored generating function. Taking absolute values in equation (6) we obtain that

$$(7) \quad \sum_{j=0}^{n-1} |\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]| |w|^j \geq |A(w)| = (1 - q) |1 - q^d| w^d + q^d \left| \left(\frac{q + (1 - q)w}{q^d + (1 - q^d)w^d} \right)^{\lfloor \frac{\ell^*}{d} \rfloor} \tilde{A}(w) \right|.$$

Ultimately, we aim to bound from below $\max_j |\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]|$ by choosing $w \in \mathbb{C}$ appropriately. To do this, we balance $|w|$ on the left hand side of equation (7) with $|(1 - q^d)w^d + q^d|$ on the right; it would be best if $|w|$ were small while $|(1 - q^d)w^d + q^d|$ were large, and so a compromise is to let w vary along an arc of the unit disc \mathbb{D} . In particular, let $\gamma_L := \{e^{i\theta} : -\pi/L \leq \theta \leq \pi/L\}$, where we assume that $L \geq 20$, a choice which will become clear later. The following lemma bounds $|(1 - q^d)w^d + q^d|$ from below while $w \in \gamma_L$. The proof is a standard calculation, deferred to Section 5.5.

LEMMA 5.2. *For $w \in \gamma_L$ we have that*

$$|(1 - q^d)w^d + q^d| \geq \exp(-2\pi^2 \cdot q^d(1 - q^d)d^2/L^2).$$

Additionally, we bound $\sup_{\gamma_L} |\tilde{A}(w)|$ from below using Lemma 5.3, whose proof is in Section 5.3.

LEMMA 5.3. *Let $0 < q < 0.7$ be a constant. There exists $\zeta \in \gamma_L$, as well as a constant $C > 0$ depending only on q , such that $|\tilde{A}(\zeta)| \geq \exp(-C \cdot dL)$.*

We are now ready to prove Theorem 1.7.

PROOF OF THEOREM 1.7. Let $\zeta \in \gamma_L$ be the point guaranteed by Lemma 5.3. Substituting ζ into equation (7), we use Lemma 5.3 and the fact that $|\zeta| = 1$ to see that

$$\begin{aligned} \sum_{j=0}^{n-1} |\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]| &\geq |A(\zeta)| = (1 - q)|(1 - q^d)\zeta^d + q^d|^{\lfloor \frac{\ell^*}{d} \rfloor} |\tilde{A}(\zeta)| \\ &\geq (1 - q)|(1 - q^d)\zeta^d + q^d|^{\lfloor \frac{\ell^*}{d} \rfloor} \exp(-C \cdot dL), \end{aligned}$$

for a constant $C > 0$ depending only on q . Using the bound $\ell^* < n$, as well as Lemma 5.2 (where we drop the factor of $1 - q^d$ in the exponent), we have that

$$\sum_{j=0}^{n-1} |\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]| \geq (1 - q) \exp(-2\pi^2 \cdot q^d nd/L^2) \exp(-C \cdot dL).$$

Setting $L = \max\{(4\pi^2 nq^d/C)^{1/3}, 20\}$ and plugging into the display above, we find that there exists an index j such that

$$(8) \quad |\mathbb{E}[b_j^1] - \mathbb{E}[b_j^2]| \geq \frac{1}{n} \exp(-C' \cdot d(nq^d)^{1/3})$$

for some constant $C' > 0$ depending only on q . Therefore, we have shown that there is some index $j = j(X^1, X^2)$ where we expect the traces corresponding to X^1 and X^2 to differ significantly.

Suppose spider X^1 goes through the deletion channel and we observe T samples, S^1, \dots, S^T where sample S^t has labels $\{u_j^t\}_{j=0}^{n-1}$. Let η denote the right hand side of equation (8). We say that a spider X^2 is a *better match* than X^1 for traces $\{S^t\}_{t \in [T]}$ if at the index $j = j(X^1, X^2)$, X^2 looks closer to the traces than X^1 ; that is, if

$$\left| \frac{1}{T} \sum_{t=1}^T u_j^t - \mathbb{E}[b_j^2] \right| \leq \left| \frac{1}{T} \sum_{t=1}^T u_j^t - \mathbb{E}[b_j^1] \right|.$$

As before, the expectation is over the random labels b^1 and b^2 . A Chernoff bound implies that if the traces $\{S^t\}_{t \in [T]}$ came from spider X^1 , then the probability that X^2 is a better match

than X^1 is at most $\exp(-T\eta^2/2)$. Repeating this for all pairs of binary labeled (n, d) -spiders, the algorithm outputs X^* , the (n, d) -spider which is a better match than all others (the best match), if such a spider exists. Otherwise, the algorithm outputs a random binary labeled (n, d) -spider.

Last, we show that the algorithm correctly reconstructs an (n, d) -spider with high probability when $d \leq \log_{1/q} n$. We bound from above the probability that the algorithm does not find that X^1 is the best match by a combination of a union bound and a Chernoff bound (as discussed above). The probabilities below are taken over the random traces $\{S^t\}_{t \in [T]}$:

$$\begin{aligned} \Pr[X^* \neq X^1] &\leq \sum_{X^2: X^2 \neq X^1} \Pr[X^2 \text{ is a better match than } X^1] \leq 2^n \cdot \exp(-T\eta^2/2) \\ &= 2^n \exp\left(-\frac{T}{2n^2} \exp(-C \cdot d(nq^d)^{1/3})\right) \end{aligned}$$

for a constant $C > 0$ depending only on q . This latter expression is at most $1/n$ if and only if

$$T \geq 2n^2(n \ln(2) + \ln(n)) \exp(Cd(nq^d)^{1/3}).$$

This holds if $T \geq \exp(cd(nq^d)^{1/3})$ for a large enough constant c depending only on q . \square

5.3. Proof of Lemma 5.3. We assume basic knowledge of subharmonic functions and harmonic measure. For background, we refer readers to any introductory complex analysis book (e.g., [2, 18]). For a more elementary (but slightly weaker) bound, see Lemma 5.9 in Section 5.5.

Let $\Omega \subset \mathbb{C}$ be a bounded, open region, and let $\partial\Omega$ denote its boundary. The *harmonic measure* of a subset $\gamma \subset \partial\Omega$ with respect to a point $w_0 \in \Omega$, denoted $\mu_\Omega^{w_0}(\gamma)$, is the probability that a Brownian motion starting at w_0 exits Ω through γ . Let $f(w)$ denote an analytic function; we will choose $f = \tilde{A}$, which is a polynomial and hence analytic. Given $|f(w_0)|$ at a point $w_0 \in \Omega$ and a condition on the growth of $|f|$ in Ω , we utilize harmonic measure to bound $|f|$ on $\partial\Omega$. Specifically, we use that $\log|f|$ satisfies the *sub-mean value property*: for all $w_0 \in \Omega$ we have that

$$(9) \quad \log|f(w_0)| \leq \int_{\partial\Omega} \log|f(w)| d\mu_\Omega^{w_0}(w).$$

As in equation (9), we will define a region of integration where the value of $\log|\tilde{A}(w)|$ is controlled along the boundary, and the boundary will contain $\gamma_L = \{e^{i\theta} : -\pi/L \leq \theta \leq \pi/L\}$. We need to separate this boundary into a few different pieces and use different techniques to upper bound $\log|\tilde{A}(w)|$ on each curve. In fact, the methods of [20] show a lower bound for $\sup_{\gamma_L} |f(w)|$ for an analytic function $f(w)$ satisfying the growth condition in Lemma 5.5 (see below), by using equation (9) and a particular choice of w_0 . We show that \tilde{A} satisfies the growth condition specified in Lemma 5.5, then borrow techniques from [20] to upper bound the right hand side of equation (9). However, we have to work more to find an appropriate point $w_0 \in \mathbb{D}$ in order to find a lower bound for the left hand side of equation (9), so that we can also show a lower bound for $\sup_{\gamma_L} |\tilde{A}(w)|$. We discuss this difficulty in Remark 3.

In what follows, open discs of radius r centered at a point z are denoted as $D_r(z)$. The unit disc, $D_1(0)$, is an exception, denoted as \mathbb{D} . Recall that $L \geq 20$, and let upper and lower case c 's with tick marks denote constants depending only on q .

PROOF OF LEMMA 5.3. First, we choose an appropriate point w_0 where we can lower bound $|\tilde{A}(w_0)|$, as stated in Lemma 5.4.

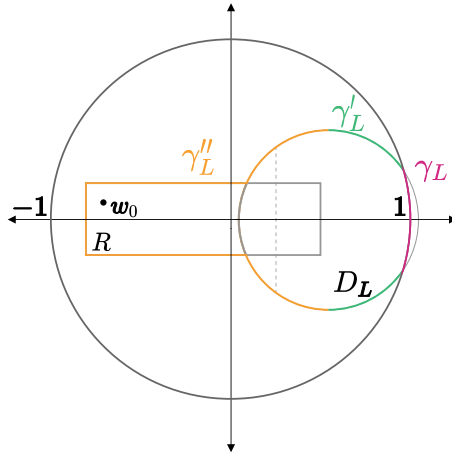


FIG. 12. The region $\Omega_L = (\mathbb{D} \cap D_L) \cup R$ with boundary $\partial\Omega_L = \gamma_L \cup \gamma'_L \cup \gamma''_L$.

LEMMA 5.4. Let $q < 0.7$ and for fixed integer $d > 0$, let $w_0 := -q$ if d is odd and $w_0 := qe^{i\pi(d-1)/d}$ if d is even. Then there exists a constant $c > 0$ depending only on q such that $|\tilde{A}(w_0)| \geq e^{-c \cdot d}$.

The calculation justifying the bound is standard and deferred to Section 5.5, but the choices of w_0 are quite careful. The choices depend on the parity of d so as to control the sign of w_0^d . With these choices we can relate $|q + (1 - q)w_0|$ to $|q^d + (1 - q^d)w_0^d|$, a motive which becomes clear in the proof of Lemma 5.4. Note that our inability to handle constant $q \geq 1/\sqrt{2}$ comes from Lemma 5.4, as detailed in Remark 3.

In the following we fix w_0 according to the specifications of Lemma 5.4. Next, we define a region Ω_L that contains w_0 , and whose boundary we integrate over; see Figure 12 for an illustration. For a translate $h_L \in \mathbb{R}^+$, let $D_L = D_{1/2}(1/2) + h_L$, where h_L is chosen so that $D_L \cap \partial\mathbb{D} = \gamma_L$. Observe that $L \geq 20$ implies $h_L \leq 1/10$. We also define a rectangle $R \subset \mathbb{D}$ that has the following properties: R contains w_0 , R has nonempty intersection with D_L , and ∂R has bounded distance from w_0 and $\partial\mathbb{D}$. As we only consider q bounded away from 1 and $d \geq 20$, we may (and will) choose R to be centered about the real axis, with height $1/5$ and with length extending from -0.8 to $1/2$. Our region of integration is then defined as $\Omega_L := (\mathbb{D} \cap D_L) \cup R$.

We partition the boundary of Ω_L into three parts by defining

$$\gamma'_L := \left\{ w \in \partial\Omega_L \setminus \gamma_L : |w| > \frac{1}{2} + h_L \right\}, \quad \gamma''_L := \left\{ w \in \partial\Omega_L : |w| \leq \frac{1}{2} + h_L \right\}.$$

We thus have that $\partial\Omega_L = \gamma_L \cup \gamma'_L \cup \gamma''_L$; see Figure 12, where the different parts of the contour are colored differently. Using the sub-mean value property, we see that

$$\begin{aligned} \log|\tilde{A}(w_0)| &\leq \int_{\partial\Omega_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \\ &= \int_{\gamma_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) + \int_{\gamma'_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) + \int_{\gamma''_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w). \end{aligned}$$

Next, we upper bound each of the integrals. Our upper bounds for γ_L and γ''_L are simple modulus length bounds. We start with the integral over γ_L . We know that the boundary $\partial\Omega_L$ has constant length, the curve γ_L has length bounded by c/L for some constant c , and w_0 is bounded away from γ_L . Therefore the probability that a Brownian motion starting at w_0 exits

Ω_L through the arc γ_L is at most C/L for some constant C . That is, $\mu_{\Omega_L}^{w_0}(\gamma_L) \leq C/L$, where we can choose C to hold for all w_0 as q and d vary. Therefore we have that

$$(10) \quad \int_{\gamma_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \leq \mu_{\Omega_L}^{w_0}(\gamma_L) \cdot \log \sup_{w \in \gamma_L} |\tilde{A}(w)| \leq \frac{C}{L} \cdot \log \sup_{w \in \gamma_L} |\tilde{A}(w)|.$$

Somewhat more work is needed to show that

$$(11) \quad \int_{\gamma'_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \leq c' \quad \text{and} \quad \int_{\gamma''_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \leq c''.$$

To prove these bounds, we use Lemma 5.5, a growth condition for the generating function.

LEMMA 5.5. *For all $w \in \mathbb{D}$ and all deletion probabilities $q \in (0, 1)$, we have $|\tilde{A}(w)| \leq \frac{1}{(1-q)(1-|w|)}$.*

The proof of Lemma 5.5 is a triangle inequality calculation that we defer to Section 5.5. We can directly apply Lemma 5.5 to bound the integral over γ''_L . By our choice of R and the fact that $h_L \leq 1/10$, we have for all $w \in \gamma''_L$ that $|w| \leq 0.83$. Applying Lemma 5.5, we see that $|\tilde{A}(w)| \leq 10/(1 - q)$ for all $w \in \gamma''_L$. Noting that harmonic measure is a probability measure and thus $\mu_{\Omega_L}^{w_0}(\gamma''_L) \leq 1$, we have that

$$\int_{\gamma''_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \leq \sup_{w \in \gamma''_L} \log|\tilde{A}(w)| \leq \log \frac{10}{1 - q}.$$

We turn now to the integral over γ'_L , where we cannot use modulus length bounds because as $|w|$ approaches 1, the factor $1/(1 - |w|)$ becomes arbitrarily large. However, we can still use Lemma 5.5 to obtain that

$$\int_{\gamma'_L} \log|\tilde{A}(w)| d\mu_{\Omega_L}^{w_0}(w) \leq \log\left(\frac{1}{1 - q}\right) + \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{\Omega_L}^{w_0}(w).$$

It remains to bound the integral on the right hand side of the display above. A Brownian motion in Ω_L starting at w_0 must hit the segment $s_L = \{w \in \Omega_L : \text{Re}(w) = \frac{1}{4} + h_L\}$ before it hits γ'_L , so

$$\int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{\Omega_L}^{w_0}(w) \leq \sup_{z \in s_L} \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{\Omega_L}^z(w).$$

Note that z and $\partial\Omega_L$ will not be arbitrarily close. Considering Brownian motions starting at z , we will upper bound the probability that it exits Ω_L through γ'_L by the probability that it exits D_L through γ'_L . As z is bounded away from $\partial\Omega_L$ and ∂D_L , the measures $\mu_{D_L}^z$ and $\mu_{\Omega_L}^z$ are equivalent, meaning they have the same null sets. Then by the Radon–Nikodym theorem, there exists a measurable function $f = d\mu_{\Omega_L}^z/d\mu_{D_L}^z$ such that for any measurable set S , $\mu_{\Omega_L}^z(S) = \int_S f(w) d\mu_{D_L}^z(w)$. From the probabilistic definition of harmonic measure, observe that there exists a constant $c > 0$ such that for any measurable $S \subset \gamma'_L$,

$$\mu_{\Omega_L}^z(S) \leq c \cdot \mu_{D_L}^z(S).$$

Since $\mu_{\Omega_L}^z/\mu_{D_L}^z$ is bounded on γ'_L , so is f up to a set of measure 0.³ The upper bound on f almost everywhere is sufficient. Then returning to our remaining integral over γ'_L , we obtain

³If f is unbounded on a set of positive measure, $A \subset \gamma'_L$, then for all C , $f \geq C$ on A . Writing $\mu_{\Omega_L}^z(A) = \int_A f d\mu_{D_L}^z \geq C \mu_{D_L}^z(A)$, we see that $\mu_{\Omega_L}^z(A)/\mu_{D_L}^z(A)$ is unbounded.

the upper bound

$$\begin{aligned} \sup_{z \in s_L} \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{\Omega_L}^z(w) &= \sup_{z \in s_L} \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) f(w) d\mu_{D_L}^z(w) \\ &\leq c \cdot \sup_{z \in s_L} \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{D_L}^z(w). \end{aligned}$$

We move to harmonic measure with respect to a disc, instead of Ω_L , so that we can switch measures to integrate with respect to angles on the disc. Specifically, we have an explicit form for the Radon–Nikodym derivative. Letting s denote the arc length measure and D_r denote a disc of radius r containing the point z , $d\mu_{D_r}^z/ds$ at a point $\zeta \in \partial D_r$ is the Poisson kernel $P(z, \zeta) = \frac{r^2 - |z|^2}{2\pi \cdot r |z - \zeta|^2}$. Note that $P(z, \zeta)$ is uniformly bounded above by a constant for all $\zeta \in \gamma'_L$, as all $z \in s_L$ are bounded away from γ'_L , since $w \in \gamma'_L$ has $\text{Re}(w) \geq 1/2$. This is a useful observation, as now we can integrate with respect to the angle on D_L between $w = (1/2 + h_L) + e^{i\theta}/2 \in \gamma'_L$ and $x = e^{i\pi/L} = 1/2 + h_L + e^{i\theta_0}/2$, while only gaining a constant factor depending on q in the upper bound. Recall x is the intersection point of ∂D_L and $\partial \mathbb{D}$ in the first quadrant. We use the following lemma to obtain a new bound. There are several proofs for it using only elementary geometry, and we include one in Section 5.5.

LEMMA 5.6. *For translate $0 < h_L \leq 1/10$ satisfying $e^{i\pi/L} \in D_L \cap \mathbb{D}$, consider the points $w = (1/2 + h_L) + e^{i\theta}/2$ with $0 \leq \theta \leq \pi/2$, and $x = e^{i\pi/L} = 1/2 + h_L + e^{i\theta_0}/2$ with $0 \leq \theta_0 \leq \theta$. Then $1 - |w| \geq \frac{1}{64}(\theta - \theta_0)^4$.*

As our region of integration is symmetric, it suffices to only show the inequality for γ'_L in the first quadrant. From the boundedness of the Poisson kernel and Lemma 5.6,

$$\begin{aligned} c \cdot \sup_{z \in s_L} \int_{\gamma'_L} \log\left(\frac{1}{1 - |w|}\right) d\mu_{D_L}^z(w) &\leq C \cdot \int_{\theta_0}^{\pi/2} \left(\log 64 + \log\left(\frac{1}{\theta - \theta_0}\right)^4\right) d\theta \\ &\leq c'' \cdot \int_{\theta_0}^{\pi/2} \log\left(\frac{1}{\theta - \theta_0}\right) d\theta \leq c'. \end{aligned}$$

Having proven the bounds on the integrals in equation (10) and equation (11), we are now ready to conclude the proof of Lemma 5.3. Combining these bounds with Lemma 5.4 and the sub-mean value property inequality, we see that

$$-cd \leq \frac{C}{L} \log \sup_{w \in \gamma_L} |\tilde{A}(w)| + c' + c'',$$

where all constants are positive and depend only on q . Rearranging, we now have the lower bound $\sup_{w \in \gamma_L} |\tilde{A}(w)| \geq \exp(-C'dL)$ for some constant C' depending only on q . As γ_L is a closed arc, there exists $\zeta \in \gamma_L$ such that $|\tilde{A}(\zeta)| = \sup_{w \in \gamma_L} |\tilde{A}(w)|$, as desired. \square

5.4. *Bounds for spiders from string trace reconstruction.* String reconstruction methods can be used as a black box for spiders. For depth $d \geq \log_{1/q} n$, this achieves the best known bound. However, for smaller depths, our algorithm is more efficient.

Large depth (n, d) -spiders.

PROOF OF PROPOSITION 1.8. With probability $(1 - q^d)^{n/d}$, a trace contains at least one non-root node from each of the n/d paths in the spider. When all paths are present, we can match paths of the trace to paths of the original spider and learn paths separately. Using

only such traces, we are faced, for each path, with a string trace reconstruction problem with censoring (see the [Appendix](#)), where the string length is d , the deletion probability is q , and the censoring probability is $\gamma = 1 - (1 - q^d)^{n/d-1}$. Lemma A.1 in the [Appendix](#) (with $\varepsilon = 1/2$) tells us that

$$T_\gamma^{\text{cens}}\left(d, \frac{1}{n^2}\right) \leq \frac{3/2}{(1 - q^d)^{n/d}} \cdot T\left(d, \frac{1}{2n^2}\right) \leq 2 \cdot T\left(d, \frac{1}{2n^2}\right),$$

where the second inequality holds (for all n large enough) because $(1 - q^d)^{n/d} \rightarrow 1$ as $n \rightarrow \infty$ when $d \geq \log_{1/q} n$. That is, if we observe $2 \cdot T(d, \frac{1}{2n^2})$ traces of the spider, then the bits along each specific path can be reconstructed with error probability at most $1/n^2$. Hence, by a union bound the bits along all paths can be reconstructed with error probability at most $(n/d) \times (1/n^2) \leq 1/n$. \square

We can extend Proposition 1.8 to the following result. We omit the proof, which follows the same outline and ideas as the proof of Proposition 1.8.

PROPOSITION 5.7. *For n large enough, $\alpha \geq 0$, and $d \geq \log_{1/q} n - \log_{1/q}(\log^{1+\alpha} n)$, an (n, d) -spider can be reconstructed with $\exp(C(\log^\alpha n)) \cdot T(d, \frac{1}{2n^2})$ traces with high probability, where C is a constant depending only on q .*

Small depth (n, d) -spiders. When $d = c \log_{1/q} n$ with constant $0 < c < 1$, the same reconstruction strategy still applies, but it does worse than our mean-based algorithm (which results in Theorem 1.7). In this regime of d , to ensure that with high probability we see even a single trace containing all n/d paths, we must take $\exp(\Omega(n^{1-c}/\log n))$ traces. It suffices to take $\exp(O(n^{1-c}/\log n)) \cdot T(d, 1/n^2) = \exp(\tilde{O}(n^{1-c}))$ traces to ensure that enough traces contain all n/d paths. However, our mean-based algorithm resulting in Theorem 1.7 does better than this, requiring only $\exp(\tilde{O}(n^{(1-c)/3}))$ traces to reconstruct.

We observe that the previous results on string trace reconstruction can also be used to derive the following proposition (in addition to Proposition 1.8 and Proposition 5.7). The consequences are twofold: (i) when $d = O(1)$, then the trace complexity of spiders is asymptotically the same as strings, and (ii) our result in Theorem 1.7 offers an improvement when $d = \omega(1)$.

PROPOSITION 5.8. *For $d < \log_{1/q} n$, we can reconstruct an (n, d) -spider with high probability by using at most $\exp((\frac{C'n}{d(1-q)^{2d}})^{1/3})$ traces, for $C' > 0$ depending on q .*

We sketch the proof. A path in the spider of depth d retains all of its nodes with probability $(1 - q)^d$. Equivalently, some node is deleted with probability $q' = 1 - (1 - q)^d$. For any trace, consider the modified channel that deletes any path entirely if it is missing at least one node. With this modification, every row of the spider behaves as if it were a string on n/d bits in a channel with deletion probability q' . Opening up the proof of Theorem 1.1, for nonconstant deletion probability $q' = 1 - (1 - q)^d$, then gives the proposition.

5.5. *Additional proofs and remarks for reconstructing spiders.*

REMARK 3 (Remark for Theorem 1.7). In the proof of Lemma 5.3, which is needed to prove Theorem 1.7, we are unable to handle general generating functions with deletion probability $1/\sqrt{2} < q < 1$. We require some anchor point, w_0 , for which we can lower bound $|\tilde{A}(w_0)|$ and a simple curve surrounding w_0 for which we can upper bound $|\tilde{A}(w)|$ along that

curve. For any fixed $|w| > 1$, for $d = \Omega(1)$ we see that $|(1 - q^d)w^d + q^d| > 1$ for sufficiently large n . This results in terms on the order of $c^{n/d}$ in our generating function, for constant $c > 1$. So our anchor point cannot lie outside of \mathbb{D} , and more specifically the surrounding curve cannot leave \mathbb{D} .

Inside the unit disc, upper bounds on $|\tilde{A}(w)|$ have a nice form due to Lemma 5.5. It seems for any fixed point in $w_0 \in \mathbb{D}$, there is a factored generating function $\tilde{A}(w)$ which is small at w_0 , $|A(w_0)| = \Theta(q^d)$. However it is not clear whether for every factored generating function $\tilde{A}(w)$ there is some $w_0 \in \mathbb{D}$, not tending to the boundary, such that $|\tilde{A}(w_0)| > c$ for some constant c depending only on q . Such arguments are common in complex analysis for families of analytic functions which are sequentially compact, but our family of generating functions does not satisfy this property.

PROOF OF LEMMA 5.1. We index the nonroot nodes of the spider according to the DFS ordering described in Section 5.1. We can uniquely write any $j \in \{0, 1, \dots, n - 1\}$ as $j = d \cdot s_j + r_j$ with $s_j \in \{0, 1, \dots, n/d - 1\}$ corresponding to a particular path of the spider and $r_j \in \{0, 1, \dots, d - 1\}$ describing where along this path node j is. Consider two nodes, $j = d \cdot s_j + r_j$ and $\ell = d \cdot s_\ell + r_\ell$, with $j \geq \ell$. After passing a through the deletion channel to get the trace b , b_ℓ comes from a_j if and only if a_j is retained, exactly r_ℓ of the first r_j nodes in the path of j are retained, and exactly s_ℓ of the first s_j paths are retained. This leads to the following generating function:

$$\begin{aligned} & \mathbb{E} \left[\sum_{\ell=0}^{n-1} b_\ell w^\ell \right] \\ &= (1 - q) \sum_{\ell=0}^{n-1} w^\ell \sum_{j=\ell}^{n-1} a_j \binom{r_j}{r_\ell} (1 - q)^{r_\ell} q^{r_j - r_\ell} \binom{s_j}{s_\ell} q^{d(s_j - s_\ell)} (1 - q^d)^{s_\ell} \mathbf{1}_{\{r_\ell \leq r_j\}} \\ &= (1 - q) \sum_{j=0}^{n-1} a_j \sum_{\ell=0}^j \binom{r_j}{r_\ell} (1 - q)^{r_\ell} q^{r_j - r_\ell} \binom{s_j}{s_\ell} q^{d(s_j - s_\ell)} (1 - q^d)^{s_\ell} w^\ell \mathbf{1}_{\{r_\ell \leq r_j\}} \\ &= (1 - q) \sum_{s_j=0}^{n/d-1} \sum_{r_j=0}^{d-1} a_{s_j d + r_j} \\ & \quad \times \sum_{s_\ell=0}^{s_j} \sum_{r_\ell=0}^{r_j} \binom{r_j}{r_\ell} (1 - q)^{r_\ell} q^{r_j - r_\ell} \binom{s_j}{s_\ell} q^{d(s_j - s_\ell)} (1 - q^d)^{s_\ell} w^{s_\ell d + r_\ell}, \end{aligned}$$

where we used linearity of expectation and interchanged the order of summation. Observing that the sums are binomial expansions we have that

$$\begin{aligned} & \mathbb{E} \left(\sum_{\ell=0}^{n-1} b_\ell w^\ell \right) \\ &= (1 - q) \sum_{s_j=0}^{n/d-1} \sum_{r_j=0}^{d-1} a_{s_j d + r_j} (q + (1 - q)w)^{r_j} (q^d + (1 - q^d)w^d)^{s_j} \\ &= (1 - q) \sum_{j=0}^{n-1} a_j (q + (1 - q)w)^{j \pmod{d}} (q^d + (1 - q^d)w^d)^{\lfloor \frac{j}{d} \rfloor}, \end{aligned}$$

which proves the claim. \square

PROOF OF LEMMA 5.2. Writing $w = \cos(\theta) + i \sin(\theta)$, we see that

$$\begin{aligned} & |(1 - q^d)w^d + q^d|^2 \\ &= |(1 - q^d)(\cos(\theta) + i \sin(\theta))^d + q^d|^2 = |(1 - q^d)(\cos(d\theta) + i \sin(d\theta)) + q^d|^2 \\ &= ((1 - q^d) \cos(d\theta) + q^d)^2 + ((1 - q^d) \sin(d\theta))^2 \\ &= (1 - q^d)^2 \cos^2(d\theta) + 2q^d(1 - q^d) \cos(d\theta) + q^{2d} + (1 - q^d)^2 \sin^2(d\theta) \\ &= (1 - q^d)^2 + 2q^d(1 - q^d) \cos(d\theta) + q^{2d} = 1 - 2q^d + 2q^{2d} + 2q^d(1 - q^d) \cos(d\theta) \\ &= 1 - 2q^d(1 - q^d)(1 - \cos(d\theta)). \end{aligned}$$

Now using the fact that $1 - \cos(y) \leq y^2/2$, as well as the inequality $1 - y \geq \exp(-4y)$ which holds for all $y \in [0, 0.9]$ (in our case indeed $q^d(1 - q^d)d^2\theta^2 \in [0, 0.9]$ for all possible parameter values), we obtain that

$$|(1 - q^d)w^d + q^d|^2 = 1 - 2q^d(1 - q^d)(1 - \cos(d\theta)) \geq \exp(-4q^d(1 - q^d)d^2\theta^2).$$

Taking a square root of the last line shows $|(1 - q^d)w^d + q^d| \geq \exp(-2q^d(1 - q^d)d^2\theta^2)$. Finally, the assumption that $w \in \gamma_L$ implies that $\theta^2 \leq \pi^2/L^2$ and the claim follows. \square

PROOF OF LEMMA 5.4. We will consider the case of even and odd d separately, starting with the cleaner case of when d is odd. Recall that we assume that $d \geq 20$ and we choose $w_0 = -q$ when d is odd and $w_0 = qe^{i \cdot \pi(d-1)/d}$ when d is even. Let $\alpha := |q + (1 - q)w_0|$ and $\beta := |q^d + (1 - q^d)w_0^d|$.

When d is odd, $\alpha = q^2$, and also $w_0^d = -q^d$, hence $\beta = q^{2d}$. When d is even, w_0 is still chosen so that $w_0^d = -q^d$ and thus $\beta = q^{2d}$, as in the case when d is odd. It is clear geometrically that $\alpha \geq q^2$, but we include the calculation as well:

$$\begin{aligned} \alpha &= \sqrt{(\operatorname{Re}(q + (1 - q)w_0))^2 + (\operatorname{Im}(q + (1 - q)w_0))^2} \\ &= \sqrt{(q + (1 - q)q \cos(\pi(d - 1)/d))^2 + ((1 - q)q \sin(\pi(d - 1)/d))^2} \\ &= \sqrt{2q^2(1 - q)(1 + \cos(\pi(d - 1)/d)) + q^4} \geq \sqrt{0 + q^4} = q^2, \end{aligned}$$

where we use that $\cos(\pi(d - 1)/d) \geq -1$. By our choice of β , we also see that $\alpha^d \geq q^{2d} = \beta$.

We are now ready to prove a lower bound on $|\tilde{A}(w_0)|$ which holds for both d even and d odd. First, recalling the definition of \tilde{A} and the fact that $w_0^d = -q^d$, we have that

$$\begin{aligned} \tilde{A}(w_0) &= (1 - q) \sum_{\ell=\ell^*}^{n-1} a_\ell (q + (1 - q)w_0)^\ell (\operatorname{mod} d) (q^d + (1 - q^d)w_0^d)^{\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor} \\ &= (1 - q) \sum_{\ell=\ell^*}^{n-1} a_\ell (q + (1 - q)w_0)^\ell (\operatorname{mod} d) q^{2d(\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor)}. \end{aligned}$$

Now recall that $|a_{\ell^*}| = 1$ and thus the first term in the sum above (corresponding to index $\ell = \ell^*$) is, in absolute value, equal to $\alpha^{\ell^* (\operatorname{mod} d)}$. Since $a_\ell \in \{-1, 0, 1\}$ for all ℓ , the rest of the sum above (adding terms corresponding to indices $\ell^* + 1 \leq \ell \leq n - 1$) is, in absolute value, at most

$$\sum_{\ell=\ell^*+1}^{n-1} \alpha^{\ell (\operatorname{mod} d)} q^{2d(\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor)} \leq \sum_{\ell=\ell^*+1}^{n-1} \alpha^{\ell (\operatorname{mod} d)} \alpha^{d(\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor)} = \alpha^{-d \lfloor \frac{\ell^*}{d} \rfloor} \sum_{\ell=\ell^*+1}^{n-1} \alpha^\ell$$

$$\leq \alpha^{-d \lfloor \frac{\ell^*}{d} \rfloor} \frac{\alpha}{1 - \alpha} \alpha^{\ell^*} = \frac{\alpha}{1 - \alpha} \alpha^{\ell^* \pmod{d}}.$$

Putting these two bounds together we obtain that

$$|\tilde{A}(w_0)| \geq (1 - q) \left(\alpha^{\ell^* \pmod{d}} - \frac{\alpha}{1 - \alpha} \alpha^{\ell^* \pmod{d}} \right) = (1 - q) \frac{1 - 2\alpha}{1 - \alpha} \alpha^{\ell^* \pmod{d}}.$$

When α is less than and bounded away from $1/2$, then this bound is at least a positive constant times $\alpha^{\ell^* \pmod{d}}$. Here, our choice of d and q becomes clear, as when $d \geq 20$ and $q \leq .7$ then $\alpha \leq .49$. Since $\ell^* \pmod{d} < d$, we have that $\alpha^{\ell^* \pmod{d}} \geq q^{2d}$ and the claim follows. \square

PROOF OF LEMMA 5.5. First, we show that $q^d + (1 - q^d)|w|^d \leq (q + (1 - q)|w|)^d$ for all $w \in \mathbb{D}$ and $q \in (0, 1)$. This is because

$$\begin{aligned} (q + (1 - q)|w|)^d &= \sum_{j=0}^d \binom{d}{j} q^j ((1 - q)|w|)^{d-j} = q^d + \sum_{j=0}^{d-1} \binom{d}{j} q^j ((1 - q)|w|)^{d-j} \\ &\geq q^d + |w|^d \sum_{j=0}^{d-1} \binom{d}{j} q^j (1 - q)^{d-j} = q^d + |w|^d (1 - q^d), \end{aligned}$$

where we used the inequality $|w|^{-j} \geq 1$ which holds when $|w| \leq 1$ and $j \geq 0$. Combining this inequality with the triangle inequality, we can show the desired upper bound for $|\tilde{A}(w)|$:

$$\begin{aligned} |\tilde{A}(w)| &\leq \sum_{\ell=\ell^*}^{n-1} |a_\ell| |q + (1 - q)w|^\ell \pmod{d} |q^d + (1 - q^d)w^d|^{\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor} \\ &\leq \sum_{\ell=\ell^*}^{n-1} (q + (1 - q)|w|)^\ell \pmod{d} (q^d + (1 - q^d)|w|^d)^{\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor} \\ &\leq \sum_{\ell=\ell^*}^{n-1} (q + (1 - q)|w|)^{\ell \pmod{d} + d(\lfloor \frac{\ell}{d} \rfloor - \lfloor \frac{\ell^*}{d} \rfloor)} \\ &= (q + (1 - q)|w|)^{-d \lfloor \frac{\ell^*}{d} \rfloor} \sum_{\ell=\ell^*}^{n-1} (q + (1 - q)|w|)^\ell \\ &\leq (q + (1 - q)|w|)^{-d \lfloor \frac{\ell^*}{d} \rfloor} \frac{(q + (1 - q)|w|)^{\ell^*}}{1 - (q + (1 - q)|w|)} \\ &\leq \frac{1}{1 - (q + (1 - q)|w|)} = \frac{1}{(1 - q)(1 - |w|)}, \end{aligned}$$

where we used that $q + (1 - q)|w| < 1$ and $\ell^* - d \lfloor \ell^*/d \rfloor \geq 0$. Note that the same upper bound holds for $|A(w)|$ as well, since $|A(w)| \leq |\tilde{A}(w)|$ for all $w \in \mathbb{D}$. \square

PROOF OF LEMMA 5.6. For the setup of this proof, it may be helpful to refer to Figure 12. $x = e^{i\pi/L}$ lies on the disc D_L , and so we can also write x as $x = 1/2 + h_L + e^{i\theta_0}/2$. On the other hand, letting $\varepsilon = 1 - |w|$ we see that $|1/2 + h_L + (1/2 + \varepsilon)e^{i\theta}| = 1$. We will assume that θ_0 and θ are in the first quadrant, and we could obtain the same result when they are both in the fourth quadrant by symmetry. Set the two moduli equal to each other:

$$(12) \quad |1/2 + h_L + 1/2e^{i\theta_0}|^2 = |1/2 + h_L + (1/2 + \varepsilon)e^{i\theta}|^2$$

Computing the left hand side of equation (12):

$$\begin{aligned} |1/2 + h_L + 1/2e^{i\theta_0}|^2 &= (1/2 + h_L + 1/2 \cos(\theta_0))^2 + (1/2 \sin(\theta_0))^2 \\ &= (1/2 + h_L)^2 + 1/4 \cos^2(\theta_0) + 1/4 \sin^2 \theta_0 + (1/2 + h_L) \cos(\theta_0) \\ &= (1/2 + h_L)^2 + 1/4 + (1/2 + h_L) \cos(\theta_0). \end{aligned}$$

Computing the right hand side of equation (12):

$$\begin{aligned} |1/2 + h_L + (1/2 + \varepsilon)e^{i\theta}|^2 &= (1/2 + h_L + (1/2 + \varepsilon) \cos(\theta))^2 + ((1/2 + \varepsilon) \sin(\theta))^2 \\ &= (1/2 + h_L)^2 + (1/2 + \varepsilon)^2 + 2(1/2 + h_L)(1/2 + \varepsilon) \cos(\theta). \end{aligned}$$

Setting the simplified terms equal we obtain that

$$1/4 + (1/2 + h_L) \cos(\theta_0) = 1/4 + \varepsilon^2 + \varepsilon + 2(1/2 + h_L)(1/2 + \varepsilon) \cos(\theta)$$

and so

$$(1/2 + h_L)(\cos(\theta_0) - \cos(\theta)) = \varepsilon^2 + \varepsilon + 2(1/2 + h_L) \cdot \varepsilon \cos(\theta) \leq 4\varepsilon.$$

Recall that $0 < h_L < 1/10$ and $0 < \theta_0 < \theta \leq \pi/2$. Then using standard identities,

$$\begin{aligned} (1/2 + h_L)(\cos(\theta_0) - \cos(\theta)) &\leq 4\varepsilon, \\ 2(1/2 + h_L) \cdot \sin((\theta - \theta_0)/2) \sin((\theta_0 + \theta)/2) &\leq 4\varepsilon. \end{aligned}$$

Using the fact that $x^2/4 \leq \sin(x/2)$ for $0 \leq x \leq \pi/2$ and $\theta \geq \theta_0$, we see that

$$\frac{(\theta - \theta_0)^2}{4} \cdot \frac{(\theta_0 + \theta)^2}{4} \leq 4\varepsilon$$

and so

$$(\theta - \theta_0)^4 \leq 64\varepsilon = 64(1 - |w|),$$

which proves the claim. \square

The following lemma and its proof are analogous to Lemma 3.1 in [34].

LEMMA 5.9. *Let $0 < q < 1/2$ be a constant. We have that*

$$\sup_{w \in \gamma_L} |\tilde{A}(w)| \geq \exp(-c \cdot dL) \cdot n^{-L},$$

where c is a constant depending only on q .

PROOF. Let $\lambda := \sup_{w \in \gamma_L} |\tilde{A}(w)|$ to simplify notation. Define the following analytic function on \mathbb{D} :

$$F(w) := \prod_{j=0}^{L-1} \tilde{A}(w \cdot e^{2\pi i j/L}).$$

Note that $F(w)$ is entire, as it is the product of polynomials. We bound $\sup_{w \in \partial \mathbb{D}} |F(w)|$ from above and below. For the upper bound, we use λ for one of the factors, and for the other $L - 1$ factors, we use the following trivial bound. For $|w| \leq 1$, the moduli of both terms in the factored generating function are at most 1, since $|w^d(1 - q^d) + q^d| \leq q^d + (1 - q^d)|w|^d \leq 1$, and so for $w \in \partial \mathbb{D}$ we have that $|\tilde{A}(w)| \leq n$. Putting these together, we obtain that $|F(w)| \leq n^{L-1}\lambda$ for all $w \in \partial \mathbb{D}$.

To obtain a lower bound for $\sup_{w \in \partial \mathbb{D}} |F(w)|$ we use the maximum principle. Observe that $|F(0)| = |\tilde{A}(0)|^L$. Since F is analytic in \mathbb{D} , by the maximum modulus principle it must achieve modulus at least $|\tilde{A}(0)|^L$ on $\partial \mathbb{D}$. Combining the upper and lower bounds on $\sup_{w \in \partial \mathbb{D}} |F(w)|$, we see that $|\tilde{A}(0)|^L \leq n^{L-1} \lambda$ and hence $\lambda \geq |\tilde{A}(0)|^L n^{-L}$.

It remains to lower bound $|\tilde{A}(0)|$. From the definition of \tilde{A} we have that

$$|\tilde{A}(0)| = (1 - q)q^{-d \lfloor \frac{\ell^*}{d} \rfloor} \left| \sum_{\ell=\ell^*}^{n-1} a_\ell q^\ell \right|.$$

Recall from the definition of ℓ^* that $|a_{\ell^*}| = 1$ and that $a_\ell \in \{-1, 0, 1\}$ for all ℓ . This implies that

$$\left| \sum_{\ell=\ell^*}^{n-1} a_\ell q^\ell \right| \geq q^{\ell^*} - \sum_{\ell=\ell^*+1}^{n-1} q^\ell = q^{\ell^*} \left(1 - \frac{q - q^{n-\ell^*}}{1 - q} \right) \geq \frac{1 - 2q}{1 - q} q^{\ell^*}.$$

Our assumption that $q < 1/2$ implies that this lower bound is positive. Putting the previous two displays together and noting that $\ell^* - d \lfloor \ell^*/d \rfloor \leq d$, we have that

$$|\tilde{A}(0)| \geq (1 - 2q)q^{\ell^* - d \lfloor \frac{\ell^*}{d} \rfloor} \geq (1 - 2q)q^d.$$

Thus we have that $\lambda \geq (1 - 2q)^L q^{dL} n^{-L}$, which proves the claim with constant $c = -\log(q - 2q^2)$. \square

6. Conclusion. We introduced the problem of tree trace reconstruction and demonstrated, for multiple classes of trees, that we can utilize the structure of trees to develop more efficient algorithms than the current state-of-the-art for string trace reconstruction. We provided new algorithms for reconstructing complete k -ary trees and spiders in two different deletion models. For sufficiently small degree or large depth, we showed that a polynomial number of traces suffice to reconstruct worst-case trees.

6.1. *Future directions.*

1. *Improved bounds.* Can our existing sample complexity bounds be improved? Our results leave open several questions for complete k -ary trees and spiders. Of particular interest are (1) the TED model for complete k -ary trees with $\omega(1) \leq k \leq c \log^2 n$ and (2) spiders with depth $d = c \log_{1/q} n$, $c < 1$; can we reconstruct with $\text{poly}(n)$ traces in these cases?

2. *General trees.* We believe our results can be extended to more general trees. In general, we do not know if the trace complexity can be bounded simply in terms of the number of nodes, the depth, and the min/max degree of the tree. What other tree structure must we take into account for tight bounds?

3. *Lower bounds.* Lower bounds have recently been proven for string trace reconstruction [10, 21]. When can analogous bounds be proven for trees? For example, is it possible to reconstruct worst-case or average-case complete binary trees with $\text{polylog}(n)$ traces?

4. *Insertions and substitutions.* We have focused on deletion channels, but insertions and substitutions are well defined and relevant for tree edit distance applications. Similarly to previous work, it would be worthwhile to understand the trace complexity for these edits.

5. *Applications.* Can insights from tree trace reconstruction be helpful in applications, for instance in computational biology? In particular, DNA sequencing and synthesis techniques are rapidly evolving, and the future statistical error correction techniques will likely be different from the ones used currently. For instance, Anavy et al. [3] recently demonstrated a new DNA storage method using composite DNA letters. Similarly, future DNA synthesis techniques may use physical constraints to enforce structure on the written bases; this could take the form of a two-dimensional array or a tree as we study.

APPENDIX

REMARK 4 (Root node is fixed). Our models implicitly enforce the property that the tree traces are connected (not forests). This is consistent with the string case, because traces are never disjoint subsequences. Moreover, in the TED model, having connected traces justifies the assumption that the root is never deleted. If the root were deleted with probability q , then the preservation of the root could be achieved by sampling $O(1/q)$ times more traces and only keeping those that are connected. In the left-propagation model, our algorithms either already reconstruct the root node as written, or could be easily altered to learn the root. This may be least obvious for spiders, but here we consider the root to be in the first path and claim that since the paths have monotonically decreasing length, our proof still holds.

It is easy to imagine other models where this assumption would not work, such as (i) allowing disconnected traces, (ii) deleting edges or subtrees, or (iii) sampling random subgraphs or graph minors. We leave such investigations as future work.

Trace reconstruction with censoring. We analyze here a variant of string trace reconstruction where each trace is independently “censored” with some probability and instead of the actual trace we receive an empty string. In other words, we have to reconstruct the original string from a random sample of the traces. Here we reduce this problem to the original string trace reconstruction problem.

Let $\mathbf{x} \in \{0, 1\}^n$ denote the original string that we aim to reconstruct. Let \emptyset denote the empty string, let $\mathcal{S}_{\geq 1} := \bigcup_{k \geq 1} \{0, 1\}^k$ denote the set of binary sequences of finite positive length, and let $\mathcal{S} := \emptyset \cup \mathcal{S}_{\geq 1}$ denote the union of this set with the empty string. Let $P_{\mathbf{x},q}$ denote the probability measure on \mathcal{S} that we obtain by passing \mathbf{x} through the deletion channel with deletion probability q . That is, if \mathbf{Y} is a random (potentially empty) string that is obtained by passing \mathbf{x} through the deletion channel, then $P_{\mathbf{x},q}(\mathbf{y}) = \Pr(\mathbf{Y} = \mathbf{y})$ for every $\mathbf{y} \in \mathcal{S}$. Now, conditionally on \mathbf{Y} , define \mathbf{Z} as follows: with probability γ , let $\mathbf{Z} = \emptyset$, and with probability $1 - \gamma$, let $\mathbf{Z} = \mathbf{Y}$. That is, $\Pr(\mathbf{Z} = \emptyset | \mathbf{Y} = \mathbf{y}) = \gamma = 1 - \Pr(\mathbf{Z} = \mathbf{y} | \mathbf{Y} = \mathbf{y})$. We call \mathbf{Z} a *censored trace* and γ the *censoring probability*. Let $P_{\mathbf{x},q,\gamma}$ denote the distribution of \mathbf{Z} on \mathcal{S} .

While in the string trace reconstruction problem we aim to reconstruct \mathbf{x} from i.i.d. traces $\mathbf{Y}_1, \dots, \mathbf{Y}_k$ from $P_{\mathbf{x},q}$, in the trace reconstruction problem with censoring we aim to reconstruct \mathbf{x} from i.i.d. censored traces $\mathbf{Z}_1, \dots, \mathbf{Z}_k$ from $P_{\mathbf{x},q,\gamma}$. Recall that $T(n, \delta)$ denotes the minimum number of traces needed to reconstruct a worst-case string on n bits with probability at least $1 - \delta$, where the dependence on the deletion probability q is left implicit.

LEMMA A.1. *Let $T_\gamma^{\text{cens}}(n, \delta)$ denote the minimum number of traces needed to reconstruct a worst-case string on n bits with probability at least $1 - \delta$ from i.i.d. censored traces from the deletion channel with deletion probability q and censoring probability γ . If $\liminf_{n \rightarrow \infty} \frac{1}{n} \log \delta > -\infty$, then for every $\varepsilon > 0$ we have that $T_\gamma^{\text{cens}}(n, \delta) \leq \frac{1+\varepsilon}{(1-q^n)(1-\gamma)} T(n, (1-\varepsilon)\delta)$ for large enough n .*

PROOF. First, note that empty strings contain no information, so reconstruction only uses nonempty traces. Next, observe that conditionally on outputting a nonempty trace, the deletion channel and the deletion channel with censoring have the same distribution. That is, if \mathbf{x} is the original string, \mathbf{Y} is a trace obtained by passing \mathbf{x} through the deletion channel, and \mathbf{Z} is a censored trace, then $\Pr(\mathbf{Y} = \mathbf{y} | \mathbf{Y} \neq \emptyset) = \Pr(\mathbf{Z} = \mathbf{y} | \mathbf{Z} \neq \emptyset)$ for every $\mathbf{y} \in \mathcal{S}_{\geq 1}$.

Suppose that we have

$$T := \frac{1 + \varepsilon}{(1 - q^n)(1 - \gamma)} T(n, (1 - \varepsilon)\delta)$$

censored traces and let T' denote the number of these censored traces that are not empty. By our first two observations we have that if $T' \geq T(n, (1 - \varepsilon)\delta)$, then we can reconstruct the original string (no matter what it was) with probability at least $1 - (1 - \varepsilon)\delta$. By construction we have that $T' \sim \text{Bin}(T, (1 - q^n)(1 - \gamma))$ and thus $\mathbb{E}T' = (1 + \varepsilon)T(n, (1 - \varepsilon)\delta)$. Hence by a Chernoff bound we have that $\Pr(T' < T(n, (1 - \varepsilon)\delta)) \leq \exp(-\frac{\varepsilon^2}{2(1+\varepsilon)^2} \mathbb{E}T') = \exp(-\frac{\varepsilon^2}{2(1+\varepsilon)} T(n, (1 - \varepsilon)\delta))$. Since $T(n, (1 - \varepsilon)\delta) = \tilde{\Omega}(n^{1.25})$ (see [21]) and $\liminf_{n \rightarrow \infty} \frac{1}{n} \log \delta > -\infty$, it follows that

$$\exp\left(-\frac{\varepsilon^2}{2(1+\varepsilon)} T(n, (1 - \varepsilon)\delta)\right) \leq \varepsilon \delta$$

for large enough n . The probability that we cannot reconstruct the original string is at most δ . \square

Note that this result is essentially optimal up to constant factors (depending on q , γ , and δ). Note also that the range of δ for which the statement holds can be relaxed (this sufficient condition was chosen for its simplicity).

Acknowledgments. We thank Nina Holden for helpful discussions relating to Lemma 5.3, and Bichlien Nguyen and Karin Strauss for pointing us to connections on branched DNA and recent work in this area. We also thank Alyshia Olsen for help designing the figures. Finally, we thank Tatiana Brailovskaya and an anonymous referee for their careful reading of the paper and their numerous helpful questions and suggestions that helped improve the paper.

An extended abstract of this paper appears in the Proceedings of the 32nd Conference on Learning Theory (COLT), 2019 [13].

Funding. The research of S.D. was supported by NSF CAREER Grant 1651861 and the David & Lucile Packard Foundation. The research of M.Z.R. was supported in part by NSF Grant DMS-1811724.

REFERENCES

- [1] ABRAHAO, B., CHERICHETTI, F., KLEINBERG, R. and PANCONESI, A. (2013). Trace complexity of network inference. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)* 491–499. ACM, New York.
- [2] AHLFORS, L. V. (1953). *Complex Analysis. An Introduction to the Theory of Analytic Functions of One Complex Variable*. McGraw-Hill Book Company, Inc., New York-Toronto-London. MR0054016
- [3] ANAVY, L., VAKNIN, I., ATAR, O., AMIT, R. and YAKHINI, Z. (2019). Data storage in DNA with fewer synthesis cycles using composite DNA letters. *Nat. Biotechnol.* **37** 1229–1236. <https://doi.org/10.1038/s41587-019-0240-x>
- [4] ANDONI, A., DASKALAKIS, C., HASSIDIM, A. and ROCH, S. (2012). Global alignment of molecular sequences via ancestral state reconstruction. *Stochastic Process. Appl.* **122** 3852–3874. MR2971717 <https://doi.org/10.1016/j.spa.2012.08.004>
- [5] BAN, F., CHEN, X., FREILICH, A., SERVEDIO, R. A. and SINHA, S. (2019). Beyond trace reconstruction: Population recovery from the deletion channel. In *Proceedings of the 60th Annual Symposium on Foundations of Computer Science (FOCS)* 745–768. IEEE, New York.
- [6] BATU, T., KANNAN, S., KHANNA, S. and MCGREGOR, A. (2004). Reconstructing strings from random traces. In *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms* 910–918. ACM, New York. MR2290981
- [7] BILLE, P. (2005). A survey on tree edit distance and related problems. *Theoret. Comput. Sci.* **337** 217–239. MR2141222 <https://doi.org/10.1016/j.tcs.2004.12.030>
- [8] BORWEIN, P. and ERDÉLYI, T. (1997). Littlewood-type problems on subarcs of the unit circle. *Indiana Univ. Math. J.* **46** 1323–1346. MR1631600 <https://doi.org/10.1512/iumj.1997.46.1435>

- [9] CEZE, L., NIVALA, J. and STRAUSS, K. (2019). Molecular digital data storage using DNA. *Nat. Rev. Genet.* **20** 456–466. <https://doi.org/10.1038/s41576-019-0125-3>
- [10] CHASE, Z. (2019). New lower bounds for trace reconstruction. Preprint. Available at <https://arxiv.org/abs/1905.03031>.
- [11] CHERAGHCHI, M., GABRYS, R., MILENKOVIC, O. and RIBEIRO, J. (2020). Coded trace reconstruction. *IEEE Trans. Inf. Theory* **66** 6084–6103. MR4173526 <https://doi.org/10.1109/TIT.2020.2996377>
- [12] CHURCH, G. M., GAO, Y. and KOSURI, S. (2012). Next-generation digital information storage in DNA. *Science* **337** 1628.
- [13] DAVIES, S., RÁČZ, M. Z. and RASHTCHIAN, C. (2019). Reconstructing trees from traces. In *Proceedings of the 32nd Conference on Learning Theory (COLT)* 961–978.
- [14] DE ANINDYA, A., O'DONNELL, R. and SERVEDIO, R. A. (2017). Optimal mean-based algorithms for trace reconstruction. In *STOC'17—Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing* 1047–1056. ACM, New York. MR3678250 <https://doi.org/10.1145/3055399.3055450>
- [15] DE ANINDYA, A., O'DONNELL, R. and SERVEDIO, R. A. (2019). Optimal mean-based algorithms for trace reconstruction. *Ann. Appl. Probab.* **29** 851–874. MR3910019 <https://doi.org/10.1214/18-AAP1394>
- [16] DUDÍK, M. and SCHULMAN, L. J. (2003). Reconstruction from subsequences. *J. Combin. Theory Ser. A* **103** 337–348. MR1996071 [https://doi.org/10.1016/S0097-3165\(03\)00103-1](https://doi.org/10.1016/S0097-3165(03)00103-1)
- [17] ERLICH, Y. and ZIELINSKI, D. (2017). DNA fountain enables a robust and efficient storage architecture. *Science* **355** 950–954. <https://doi.org/10.1126/science.aaj2038>
- [18] GARNETT, J. B. and MARSHALL, D. E. (2005). *Harmonic Measure. New Mathematical Monographs* **2**. Cambridge Univ. Press, Cambridge. MR2150803 <https://doi.org/10.1017/CBO9780511546617>
- [19] GOLDMAN, N., BERTONE, P., CHEN, S., DESSIMOZ, C., LEPROUST, E. M., SIPOS, B. and BIRNEY, E. (2013). Towards practical, high-capacity, low-maintenance information storage in synthesized DNA. *Nature* **494** 77–80.
- [20] HARTUNG, L., HOLDEN, N. and PERES, Y. (2018). Trace reconstruction with varying deletion probabilities. In *2018 Proceedings of the Fifteenth Workshop on Analytic Algorithmics and Combinatorics (ANALCO)* 54–61. SIAM, Philadelphia, PA. MR3773635 <https://doi.org/10.1137/1.9781611975062.6>
- [21] HOLDEN, N. and LYONS, R. (2020). Lower bounds for trace reconstruction. *Ann. Appl. Probab.* **30** 503–525. MR4108114 <https://doi.org/10.1214/19-AAP1506>
- [22] HOLDEN, N., PEMANTLE, R. and PERES, Y. (2018). Subpolynomial trace reconstruction for random strings and arbitrary deletion probability. In *Proceedings of the 31st Conference on Learning Theory (COLT)* 1799–1840.
- [23] HOLENSTEIN, T., MITZENMACHER, M., PANIGRAHY, R. and WIEDER, U. (2008). Trace reconstruction with constant deletion probability and related results. In *Proceedings of the Nineteenth Annual ACM-SIAM Symposium on Discrete Algorithms* 389–398. ACM, New York. MR2487606
- [24] KARAU, P. and TABARD-COSSA, V. (2018). Capture and translocation characteristics of short branched DNA labels in solid-state nanopores. *ACS Sens* **3** 1308–1315. <https://doi.org/10.1021/acssensors.8b00165>
- [25] KELLY, P. J. (1957). A congruence theorem for trees. *Pacific J. Math.* **7** 961–968. MR0087949
- [26] KRASIKOV, I. and RODITTY, Y. (1997). On a reconstruction problem for sequences. *J. Combin. Theory Ser. A* **77** 344–348. MR1429086 <https://doi.org/10.1006/jcta.1997.2732>
- [27] KRISHNAMURTHY, A., MAZUMDAR, A., MCGREGOR, A. and PAL, S. (2019). Trace reconstruction: Generalized and parameterized. In *27th Annual European Symposium on Algorithms. LIPIcs. Leibniz Int. Proc. Inform.* **144** Art. No. 68, 25. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern. MR4007907
- [28] LAURI, J. and SCAPELLATO, R. (2016). *Topics in Graph Automorphisms and Reconstruction*, 2nd ed. London Mathematical Society Lecture Note Series **432**. Cambridge Univ. Press, Cambridge. MR3496604 <https://doi.org/10.1017/CBO9781316669846>
- [29] LEVENSHTAIN, V. I. (2001). Efficient reconstruction of sequences from their subsequences or supersequences. *J. Combin. Theory Ser. A* **93** 310–332. MR1805300 <https://doi.org/10.1006/jcta.2000.3081>
- [30] MARANZATTO, T. J. (2020). Tree trace reconstruction: Some results. Thesis, New College of Florida.
- [31] MCGREGOR, A., PRICE, E. and VOROTNIKOVA, S. (2014). Trace reconstruction revisited. In *Algorithms—ESA 2014. Lecture Notes in Computer Science* **8737** 689–700. Springer, Heidelberg. MR3253172 https://doi.org/10.1007/978-3-662-44777-2_57
- [32] MITZENMACHER, M. (2009). A survey of results for deletion channels and related synchronization channels. *Probab. Surv.* **6** 1–33. MR2525669 <https://doi.org/10.1214/08-PS141>
- [33] MOSSEL, E. and ROSS, N. (2019). Shotgun assembly of labeled graphs. *IEEE Trans. Netw. Sci. Eng.* **6** 145–157. MR3969756 <https://doi.org/10.1109/TNSE.2017.2776913>

- [34] NAZAROV, F. and PERES, Y. (2017). Trace reconstruction with $\exp(O(n^{1/3}))$ samples. In *STOC'17—Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing* 1042–1046. ACM, New York. [MR3678249](#)
- [35] ORGANICK, L., ANG, S. D., CHEN, Y.-J., LOPEZ, R., YEKHANIN, S., MAKARYCHEV, K., RACZ, M. Z., KAMATH, G., GOPALAN, P. et al. (2018). Random access in large-scale DNA data storage. *Nat. Biotechnol.* **36** 242–248.
- [36] ULAM, S. M. (1960). *A Collection of Mathematical Problems*. Interscience Tracts in Pure and Applied Mathematics, No. 8. Interscience Publishers, New York. [MR0120127](#)
- [37] VISWANATHAN, K. and SWAMINATHAN, R. (2008). Improved string reconstruction over insertion-deletion channels. In *Proceedings of the Nineteenth Annual ACM-SIAM Symposium on Discrete Algorithms* 399–408. ACM, New York. [MR2487607](#)
- [38] YAZDI, S. H. T., GABRYS, R. and MILENKOVIC, O. (2017). Portable and error-free DNA-based data storage. *Sci. Rep.* **7** 5011.
- [39] YAZDI, S. H. T., KIAH, H. M., GARCIA-RUIZ, E., MA, J., ZHAO, H. and MILENKOVIC, O. (2015). DNA-based storage: Trends and methods. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications* **1** 230–248.
- [40] ZHANG, K. and SHASHA, D. (1989). Simple fast algorithms for the editing distance between trees and related problems. *SIAM J. Comput.* **18** 1245–1262. [MR1025472](#) <https://doi.org/10.1137/0218082>