

HHS Public Access

Author manuscript

IEEE Trans Inf Theory. Author manuscript; available in PMC 2019 September 19.

Published in final edited form as:

IEEE Trans Inf Theory. 2018 September; 64(9): 6070–6080. doi:10.1109/TIT.2018.2851224.

Lossless Compression of Binary Trees with Correlated Vertex Names

Abram Magner,

NSF Center for the Science of Information, Purdue University, West Lafayette, IN 47907.

Krzysztof Turowski,

NSF Center for the Science of Information, Purdue University, West Lafayette, IN 47907.

Department of Computer Science, Purdue University, IN 47907, USA (spa@cs.purdue.edu)

Faculty of Electronics, Telecommunications and Informatics, Gda sk University of Technology, Poland.

Wojciech Szpankowski [Fellow, IEEE]

NSF Center for the Science of Information, Purdue University, West Lafayette, IN 47907.

Abram Magner: abram10@gmail.com; Krzysztof Turowski: krzysztof.szymon.turowski@gmail.com; Wojciech Szpankowski: spa@cs.purdue.edu

Abstract

Compression schemes for advanced data structures have become a central modern challenge. Information theory has traditionally dealt with conventional data such as text, images, or video. In contrast, most data available today is multitype and context-dependent. To meet this challenge, we have recently initiated a systematic study of advanced data structures such as unlabeled graphs [8]. In this paper, we continue this program by considering trees with statistically correlated vertex names. Trees come in many forms, but here we deal with binary plane trees (where order of subtrees matters) and their non-plane version (where order of subtrees doesn't matter). Furthermore, we assume that each name is generated by a known memoryless source (horizontal independence), but a symbol of a vertex name depends in a Markovian sense on the corresponding symbol of the parent vertex name (vertical Markovian dependency). Such a model is closely connected to models of phylogenetic trees. While in general the problem of multimodal compression and associated analysis can be extremely complicated, we find that in this natural setting, both the entropy analysis and optimal compression are analytically tractable. We evaluate the entropy for both types of trees. For the plane case, with or without vertex names, we find that a simple two-stage compression scheme is both efficient and optimal. We then present efficient and optimal compression algorithms for the more complicated non-plane case.

I. INTRODUCTION

Over the last decade, repositories of various data have grown enormously. Most of the available data is no longer in conventional form, such as text or images. Instead, biological data (e.g., gene expression data, protein interaction networks, phylogenetic trees), topographical maps (containing various information about temperature, pressure, etc.), medical data (cerebral scans, mammogram, etc.), and social network data archives are in the

form of *multimodal* data structures, that is, multitype and context dependent structures (see Figure 1). For efficient compression of such data structures, one must take into account not only several different types of information, but also the statistical dependence between the general data labels and the structures themselves. In Figure 1 we show an example of such a multimodal structure representing an annotated protein interaction network in which graphs, trees, DAGs, and text are involved.

This paper is a step in the direction of the larger goal of a unified framework for compression of multimodal graph and tree structures. We focus on compression of trees with structure-correlated vertex "names" (strings). There are myriad examples of trees with correlated names. As an example, in Figure 2 we present a Linnaean taxonomy of *hominoid*. As is easy to see, names lower in the tree (children) are (generally) variations of the name at the root of a subtree. Closer to the models that we will consider, there are binary tree models of speciation in phylogenetics, in which each vertex (representing an organism or a species) is associated with a genome (a DNA string); a parent-child relationship in this tree indicates a biological parent-child relationship (either at the level of species or individual organisms, depending on the application context). Naturally, the genome of a child is related to that of its parent; furthermore, the sequence of genomes in a given path down the tree has a Markov property: a child's genome is related to that of its ancestors only through that of its parent.

To capture this intuitive behavior, we first formalize two natural models for random binary trees with correlated names and prove that they are equivalent (see Theorem 1 and Corolary 1).

We focus on two variations: plane-oriented and non-plane- oriented trees [3]. In planeoriented trees, the order of subtrees matters, while in non-plane-oriented trees (e.g., representing a phylogenetic tree [20]) all orientations of subtrees are equivalent (see Figure 3). For the plane-oriented case, we give an exact formula for the entropy of the labeled tree (see Theorem 2) and find that a simple two-stage compression algorithm (one stage to encode the tree structure, and the second to encode the names) based on arithmetic coding is both efficient and optimal to within two bits of the entropy. The observation that the natural two-stage scheme is optimal (so that the additional complication induced by the multimodality of this setting is only *apparent*) is of interest in light of the connection to phylogenetic trees and the fact that, in general (e.g., for various graph models with names), the problem of multimodal compression can be analytically and algorithmically intractable.

In the more complicated non-plane case, we focus on unlabeled trees and first derive the entropy rate (see Theorem 3). Then we propose two compression algorithms: a simple one of time complexity O(n) (where *n* is the number of leaves) that achieves compression within a multiplicative factor of 1% of the entropy, and then an optimal algorithm of time complexity $O(n^2)$ that is within two bits from the entropy. We note that non-plane trees are more challenging to study due to the need to account for the typically large number of symmetries. In particular, the entropy for non-plane trees depends on what we call the *tree Rényi entropy* of order 2 discussed in the remark below Theorem 3.

Regarding prior work, literature on tree and graph compression is quite scarce. For unlabeled graphs there are some recent information-theoretic results, including [7], [8] (see also [10]) and [9]. In 1990, Naor [7] proposed an efficiently computable representation for unlabeled graphs (solving Turán's [11] open question) that is optimal up to the first two leading terms of the entropy when all unlabeled graphs are equally likely. Naor's result is asymptotically a special case of recent work of Choi and Szpankowski [8], who extended Naor's result to general Erd s-Rényi graphs. In particular, in [8] the entropy and an optimal compression algorithm (up to two leading terms of the entropy) for Erd s-Rényi graph structures were presented. Furthermore, in [9] an automata approach was used to design an optimal graph compression including grammar-based compression for some data structures. Peshkin [23] proposed an algorithm for a graphical extension of the one-dimensional SEQUITUR compression method. However, SEQUITUR is known not to be asymptotically optimal. For binary plane-oriented trees rigorous information-theoretic results were obtained in [13], complemented by a universal grammar-based lossless coding scheme [14].

However, for structures with vertex names, which is the topic of this paper, there have been almost no attempts at theoretical analyses, with the notable exception of [12] for sparse Erd s-Rényi graphs. The only significant algorithmic results have been based on heuristics, exploiting the well- known properties of special kinds of graphs: for example in the case of Web graphs, their low degree and clustering [24] (see also [25], [27]). Similarly, there are some algorithms with good practical compression rate for phylogenetic trees (see [26]); however, they too lack any theoretical guarantees of their performance.

The tree models (without vertex names) that we consider in this paper are variations of the *Yule* distribution [17], which commonly arises in mathematical phylogenetics. Various aspects of these tree models have been studied in the past (see, e.g., [16]), but the addition of vertex names and the consideration of information-theoretic questions about the resulting models seems to be a novel aspect of the present work.

The rest of the paper is structured as follows: in Section II, we formulate the models and present the main results, including entropy formulas/bounds and performances of the compression algorithms. In Section III, we present the derivations of the entropy results. In Section IV, we present the compression algorithms and their analyses. Finally, in Section V, we conclude with future directions.

II. MAIN THEORETICAL RESULTS

In this section we introduce the concepts of binary plane and non-plane trees together with the notion of locally correlated names associated with their vertices. Then we define two models for these types of trees with correlated names and show the equivalence of these two models.

A. Basic definitions and notation

We call a rooted tree a *plane tree* when we distinguish left- to-right-order of the children of the nodes in the embedding of a tree on a plane (see [3]). To avoid confusion, we call a tree

with no fixed ordering of its subtrees a *non-plane* tree (also known in the literature as *Otter trees* [20]). In a non-plane tree any orientation of subtrees is equivalent.

Let *T* be the set of all binary rooted plane trees having finitely many vertices and, for each positive integer *n*, let T_n be the subset of *T* consisting of all trees with exactly *n* leaves. Similarly, let *S* and *S_n* be the set of all binary rooted non-plane trees with finitely many vertices and exactly *n* leaves, respectively.

We can also augment our trees with vertex names - given the alphabet A, names are simply words from A^m for some integer m 1. Let LT_n and LS_n be the set of all binary rooted plane and non-plane trees with names, respectively, having exactly n leaves with each vertex assigned a name - a word from A^m . In this paper we consider the case where names are correlated with the structure as discussed below.

Formally, a rooted plane binary tree $t \in T_n$ can be uniquely described by a set of triples (v_i, v_j, v_k) , consisting of a parent vertex and its left and right children. Given such a set of triples defining a valid rooted plane binary tree, we can standardize it to a tree defined on the vertex set $[2n - 1] = \{1, ..., 2n - 1\}$ by replacing vertex v in each triple by the depth-first search index of v. We then consider two trees to be the same if they have the same standard representation. See the example in Figure 4. Furthermore, a tree with names $It \in LT_n$ can be uniquely described as a pair (t, f), where t is a tree, described as above, and f is a function from the vertices of t to the words in A^m . We can also describe an element of LT_n as a set of pairs (v_i, I_i) for all i = 1, 2, ..., 2n - 1 representing the vertices and their names. If we identify the vertices (v_i) with the integers 1, 2,..., 2n - 1, f can be described as a sequence of 2n - 1 words from A^m .

Tree Notation.—Here we give notation regarding trees that will be used throughout the paper. Let *t* be a binary plane tree, and consider a vertex *v* in *t*.

We denote by $\lambda(t, v)$ and p(t, v) the left and right child of v in t, respectively (i.e., these are vertices in t). By t(v), we shall mean the subtree of t rooted at v. We denote by t^L and t^R the left and right subtree of t, respectively. We denote by (t) the number of leaves in the tree t. Finally, we denote by *root*(t) the root node of t.

B. The model

We now present a model for generating plane trees with structure-correlated names. This model will be such that individual names are tuples of independent and identically distributed symbols, (a property which we shall call *horizontal* independence), but the letters of names of children depend on the name of the parent (vertical Markovian dependence).

Our main model MT_1 is defined as follows: given the number *n* of leaves in the tree, the length of the names *m*, the alphabet *A* (of size |A|) and the transition probability matrix *P* of size $|A| \times |A|$ (representing an ergodic Markov chain) with its stationary distribution π (i.e. $\pi P = \pi$), we define a random variable LT_n as a result of the following process: starting from a single node with a randomly generated name of length *m* by a memoryless source with distribution π , we repeat the following steps until a tree has exactly *n* leaves:

- pick uniformly at random a leaf v in the tree generated in the previous step,
- append two children v_L and v_R to v_r ,
- generate correlated names for v_L and v_R by taking each letter from v and generating new letters according to P. for every letter of the parent we pick the corresponding row of matrix P and generate randomly the respective letters for v_L and v_R .

Alternatively, LT_n is equivalent to an ordered pair of random variables (T_n, F_n) , where T_n is a random variable supported on T_n and F_n is a random variable supported on sequences of words from A^m of length 2n - 1.

Our second model, MT_2 (also known as the binary search tree model), is ubiquitous in the computer science literature, arising for example in the context of binary search trees formed by inserting a random permutation of [n - 1] into a binary search tree [5]. Under this model we generate a random tree T_n as follows: *t* is equal to the unique tree in T_1 and we associate a number *n* with its single vertex. Then, in each recursive step, let $v_1, v_2, ..., v_k$ be the leaves of *t*, and let integers $n_1, n_2, ..., n_k$ be the values assigned to these leaves, respectively. For each leaf v_i with value $n_i > 1$, randomly select an integer s_i from the set $\{1, ..., n_i - 1\}$ with probability $\frac{1}{n_i - 1}$ (independently of all other such leaves), and then grow two edges from v_i

with left edge terminating at a leaf of the extended tree with value s_i and right edge terminating at a leaf of the extended tree with value $n_i - s_i$. The extended tree is the result of the current recursive step. Clearly, the recursion terminates with a binary tree having exactly n leaves, in which each leaf has assigned value 1; this tree is T_n .

The assignment of names to such trees, given length m, alphabet A, transition matrix P and its stationary distribution π , proceeds exactly as in the MT_1 model: we first generate a random word from A^m by a memoryless source with distribution π and assign it to a root, then generate correlated names for children, according to the names of the parents and to the probabilities in P. Throughout the paper, we will write P(y|x), for symbols $x, y \in A$, as the probability of transitioning from x to y according to the Markov chain P. We will also abuse notation and write P(x) as the probability assigned to x by π .

Recall that (*t*) is the number of leaves of a tree *t*, and (*t*(*v*)) denotes the number of leaves of a tree *t* rooted at *v*. It is easy to see [13] that under the model MT_2 , $P(T_1 = t_1) = 1$ (where t_1 is the unique tree in T_1) and

$$P(T_n = t) = \frac{1}{n-1} P\left(T_{\Delta(t^L)} = t^L\right) P\left(T_{\Delta(t^R)} = t^R\right),$$

which leads us to the formula

$$P(T_n = t) = \prod_{v \in V(t)} (\Delta(t(v)) - 1)^{-1} \quad (1)$$

where $\mathring{V}(t)$ is the set of the internal vertices of *t*. It turns out that the probability distribution in the MT_1 model is exactly the same, as we prove below.

Theorem 1. Under the model MT_1 it holds that

$$P(T_n = t) = \prod_{v \in V(t)} (\Delta(t(v)) - 1)^{-1} \quad (2)$$

where V(t) is the set of the internal vertices of t.

Proof: We follow [17] and use the concept of *labeled histories:* a pair (t, ξ), where $t \in T_n$ (generated by the model MT_1) and ξ is a permutation of the natural numbers from 1 to n - 1, assigned to the *internal* n - 1 vertices of t, such that every path from the root forms an ascending sequence. If we think of the internal nodes of t as being associated with their DFS numbers, then the permutation ξ is a function mapping each internal node of t to the time at which its two children were added to it. Thus, for instance, the root always receives the number 1, meaning that $\xi(1) = 1$. Clearly, for each t, each labeled history (t, ξ) corresponds to exactly one sequence of generation, as it defines uniquely the order of the leaves, which are picked during consecutive stages of the algorithm. Moreover, the probability of each feasible labeled history (t, ξ) is equal to $\frac{1}{(n-1)!}$ since it involves choosing one leaf from k available at the kth stage of the algorithm for k = 1, ..., n-1. Therefore, denoting $q(t) = |\{\xi: (t, \xi) is a labelled history \}|$, we find

$$P(T_n = t) = \frac{q(t)}{(n-1)!}$$

Note that if $\Delta(t(v)) = k$ for any vertex v of t, then we know that the sequence of node choices corresponding to (t, ξ) must contain exactly k - 1 internal vertices from the subtree of v, and that v is the first of them. Moreover, for the subsequence in the sequence corresponding to (t, ξ) of a subtree t(v) rooted at vertex v, the sequences of $\Delta(t(v)^L) - 1$ vertices from its left subtree and of $\Delta(t(v)^R) - 1$ vertices from its right subtree are interleaved in any order. Thus we arrive at the following recurrence for q(t):

$$q(t) = \begin{pmatrix} \Delta(t) - 2\\ \Delta(t^L) - 1 \end{pmatrix} q(t^L)q(t^R)$$

$$\frac{(\Delta(t)-2)!}{\left(\Delta\left(t^{L}\right)-1\right)!\left(\Delta\left(t^{R}\right)-1\right)!}q(t^{L})q(t^{R})$$

This recurrence can be solved by observing that each internal vertex appears exactly once in the numerator, and that each internal vertex not equal to the root *r* appears exactly once in the denominator. Hence

$$q(t) = \frac{\prod\limits_{v \in \hat{v}(t)} (\Delta(t(v)) - 2)!}{\prod\limits_{v \in \hat{v}(t)/\{r\}} (\Delta(t(v)) - 1)!}$$

Since $\Delta(t(r)) = \Delta(t) = n$, we thus obtain

$$q(t) = \frac{\prod_{v \in V(t)}^{n} (\Delta(t(v)) - 2)!}{\prod_{v \in V(t) \setminus \{r\}}^{n} (\Delta(t(v)) - 1)!} = \frac{(n-2)!}{\prod_{v \in V(t) \setminus \{r\}}^{n} (\Delta(t(v)) - 1)}$$

leading finally to

$$P(Tn = t) = \frac{q(t)}{(n-1)!} = \frac{1}{\prod_{v \in V(t)} (\Delta(t(v)) - 1)}$$

which completes the proof. \blacksquare

Clearly, both models are equivalent, since the underlying binary plane trees have exactly the same probability distributions:

Corollary 1. *The models* MT_1 *and* MT_2 *are equivalent in the sense that they lead to the same probability distribution on trees.*

We now can extend the above equivalence to non-plane trees. We define models MS_1 and MS_2 for non-plane trees: just generate the tree according to MT_1 (respectively, MT_2) and then treat the resulting plane tree T_n as a non-plane one S_n . Since MT_1 and MT_2 are equivalent, so are MS_1 and MS_2 . In graph terminology, MS_1 and MS_2 are unlabeled version of their MT counterparts [8].

III. ENTROPY EVALUATIONS

In this section we estimate the entropy for the plane-oriented trees with names and that of non-plane trees without names. We shall see that the latter problem is mathematically more challenging.

A. The entropy of the plane trees with names

Recall that $\Delta(LT_n^L)$ is a random variable corresponding to the number of leaves in the left subtree of LT_n . From the previous section, we know that $P(\Delta(LT_n^L) = i) = \frac{1}{n-1}$ Let also r_n denote the root of LT_n , for each n.

First, by the chain rule for entropy, we can write

$$H(LT_n) = H(T_n) + H(F_n | T_n) \quad (3)$$

The second term, which is the entropy of the names given the tree T_n , is easy to estimate as a result of the stationarity of the Markov process generating the vertex names. In particular,

$$H(F_n | T_n) = 2mh(P)(n - 1) + mh(\pi)$$

where $h(P) = -\sum_{a \in A} \pi(a) \sum_{b \in A} P(b|a) \log P(b|a)$ is the entropy of the Markov process with transition matrix *P* and stationary distribution π . It is also the entropy of a child letter given its parent letter. The second term in the above expression is simply the entropy of the name of the root. We note an important subtlety here: while the label function F_n is *statistically* dependent on the tree T_n the conditional entropy $H(F_n|T_n = t_n)$, viewed as a function of the tree structure t_n is constant.

It remains to estimate $H(T_n)$ in (3). In fact, this was already computed in [13]. This leads to our first main result.

Theorem 2. The entropy of a binary tree with names of fixed length m, generated according to the model MT_1 , is given by

$$H(LT_n) = \log(n-1) + 2n \sum_{k=2}^{n-1} \frac{\log_2(k-1)}{k(k+1)} + 2mh(P)(n-1) + mh(\pi)$$
(4)

where $h(\pi) = -\sum_{a \in A} \pi(a) log \pi(a)$

Remark. Note that because F_n and T_n are statistically dependent, $H(LT_n)$ is *strictly less than* $H(F_n) + H(T_n)$, the fundamental lower bound on the expected code length if the names and the tree structure are compressed separately. This is to be expected: if one is given the multiset of names without any additional information about the tree structure, then the compression of the names either duplicates information from the tree (in the case where the tree can be recovered from the names) or is suboptimal in light of the fact that the dependence structure of the names cannot be estimated. In either case, separate compression is suboptimal.

Remark. We know that (see also [13])

$$2\sum_{k=2}^{n} \frac{\log_2(k-1)}{k(k+1)} \approx 1.736$$

for large *n* (we took $n = 10^6$). The above series converges slowly, with an error term of order $O(\log n/n)$. Therefore one would prefer a more explicit formula, which we offer next for large *n*. We approximate the sum using the Euler-Maclaurin formula [4] by the integral

$$\sum_{k=1}^{n-2} \frac{\log k}{(k+1)(k+2)} \sim \int_{1}^{n-2} \frac{\log(x)}{(x+1)(x+2)} dx = -Li_2 \left(1 - \frac{n}{2}\right) + Li_2 (2-n) + \log\left(2 - \frac{2}{n}\right) \cdot \log(n-2) + Li_2 \left(-\frac{1}{2}\right) + \frac{\pi^2}{12}$$

where

$$Li_2(x) = \int_1^x \frac{\log t}{1-t} dt$$

is the *dilogarithmic integral*. For large x the following holds [18], [19]

$$Li_{2}(x) = -\frac{1}{2}log^{2}x - \frac{\pi^{2}}{6} + O(logx/x).$$

In fact, to get a better approximation of the sum we need two extra terms in the Euler-Maclaurin formula which leads to the following approximation

$$\sum_{k=1}^{n-2} \frac{\log k}{(k+1)(k+2)} \approx Li_2(\frac{3}{2}) + \frac{1}{12}\pi^2 + \frac{1}{2}\log^2(2) - \frac{1}{72} + \frac{23}{12960} = 0.868..$$

which matches the first three digits of the sum.

Remark Note that the total entropy is proportional to *n* and *m*. In a typical setting m = O(1) or $\Theta(\log n)$ (which is certainly needed if we want all the labels to be distinct), so $H(LT_n)$ would be O(n) or $O(n \log n)$, respectively.

An alternative way to prove Theorem 2 proceeds by writing $H(LT_n) = H(LT_n|F_n(r_n)) + H(F_n(r_n))$, and the first term turns out to satisfy a recurrence which will also appear in a more fundamental way in several other places in the paper. Therefore, we state the following technical lemma, yielding the solution to this recurrence.

Lemma 1. *The recurrence* $x_1 = 0$,

$$x_n = a_n + \frac{2}{n-1} \sum_{k=1}^{n-1} x_k, n \ge 2 \quad (5)$$

has the following solution for n 2:

$$x_n = a_n + n \sum_{k=2}^{n-1} \frac{2a_k}{k(k+1)}$$
(6)

Proof: By comparing the equations for x_n and x_{n+1} for any n-2 we obtain

$$\frac{x_{n+1}}{n+1} = \frac{x_n}{n} + \frac{na_{n+1} - (n-1)a_n}{n(n+1)}$$

Substituting $y_n = -\frac{x_n}{n}$ and $b_n = \frac{na_{n+1} - (n-1)a_n}{n(n+1)}$ we get

$$y_{n+1} = y_n + b_n = y_2 + \sum_{k=2}^n b_k$$

$$y_n = y_2 + \sum_{k=2}^{n-1} b_k = y_2 + \sum_{k=2}^{n-1} \left(\frac{a_k + 1}{k+1} - \frac{a_k}{k} + \frac{2a_k}{k(k+1)} \right) = y_2 + \frac{a_n}{n} - \frac{a_2}{2} + \sum_{k=2}^{n-1} \frac{2a_k}{k(k+1)} = \frac{a_n}{n} + \sum_{k=2}^{n-1} \frac{2a_k}{k(k+1)}$$

and finally, after multiplying both sides by *n*, we obtain the desired result.

Note that Theorem 2 follows directly from Lemma 1.

B. The entropy of the non-plane trees

Now we turn our attention to the non-plane trees and the case when m = 0 or, equivalently, S_n instead of LS_n . The case of labeled non-plane trees may be easily derived from this, because of the following relation:

$$H(LS_n) = H(S_n) + H(F_n | S_n)$$

where

$$H(F_n | S_n) = H(F_n | T_n) = 2mh(P)(n-1) + mh(\pi)$$

just as in plane trees. Thus,

$$H(LS_n) = H(T_n) - H(T_n | S_n) + H(F_n | T_n)$$

may easily be turned into an explicit expression.

We thus only consider only the unlabeled case in what follows. Let S_n be the random variable with probability distribution given by the MS_2 (or, equivalently, MS_1) model. For any $s \in S$ and $t \in T$ let $t \sim s$ mean that the plane tree t is isomorphic to the non-plane tree s. Furthermore, we use the following notation: $[s] = \{t \in T: t \sim s\}$. For any $t_1; t_2 \in T_n$ such that $t_1 \sim s$ and $t_2 \sim s$ it holds that $P(T_n - t_1) = P(T_n = t_2)$, since there is also an isomorphism between them [17]. By definition, s corresponds to |[s]| isomorphic plane trees, so for any $t \sim s$ it holds that

$$P(S_n = s) = |[s]| P(T_n = t)$$
 (7)

$$P(T_n = t | S_n = s) = \frac{1}{|[s]|}.$$
 (8)

Let us now introduce two functions: X(t) and Y(t), which are equal to the number of internal vertices of $t \in T$ with unbalanced subtrees (unbalanced meaning that the number of leaves in its left and right subtree are not equal) and the number of internal vertices with balanced, but non isomorphic subtrees, respectively. Similarly, let X(s) and Y(s) denote the number of such vertices for $s \in S$. Clearly, for any $t \in T$, $s \in S$ such that $t \sim s$ it is straightforward that X (t) = X(s) and Y(t) = Y(s). Moreover, observe that any structure $s \in S$ has exactly

$$[s] = 2^{X(s) + Y(s)}$$

distinct plane orderings, since each vertex with non isomorphic subtrees can have either one of two different orderings of the subtrees, whereas when both subtrees are isomorphic we have only one ordering.

Now we conclude that

$$H(T_n | S_n)$$

$$= -\sum_{t \in T_n, s \in S_n} P(T_n = t, S_n = s) log P(T_n = t | S_n = s) = \sum_{s \in S_n, t \sim s} P(S_n = s) log |[s]|$$

$$= \sum_{s \in S_n, t \sim s} P(S_n = s) (X(s) + Y(s)) = \sum_{t \in T_n} P(T_n = t) (X(t) + Y(t)) = EX_n + EY_n$$

where we write $X_n := X(T_n)$ and $Y_n := Y(T_n)$. We thus need to evaluate $EX_n + EY_n$. It turns out to be easiest to do this by determining the expected value of a random variable Z_n representing the number of internal vertices with isomorphic child subtrees in a tree of size *n*. It satisfes $Z_n = n - 1 - (X_n + Y_n)$. However, first we determine $E[X_n]$, which will be useful later, in the proof of Theorem 4.

The stochastic recurrence for X_n can be expressed as follows: $X_1 = 0$ and

$$X_{n} = X_{U_{n-1}} + X_{n-Un-1} + I\left(U_{n-1} \neq \frac{n}{2}\right)$$

for n = 2, where U_n is the variable on $\{1, 2, ..., n\}$ with uniform distribution, and the indicator (denoted by $I(\bullet)$) contributes whenever the left and right subtrees of the root are of unequal size. This leads us to the following formula for EX_n

$$\begin{split} & EX_n = \frac{1}{n-1} \sum_{k=1}^{n-1} E \Big(X_n / U_{n-1} = k \Big) = \frac{1}{n-1} \sum_{k=1}^{n-1} \Big(EX_k + EX_{n-k} + E \Big(I \Big(U_{n-1} \neq \frac{n}{2} \Big) U_{n-1} = k \Big) \Big) \\ & = EI \Big(U_{n-1} \neq \frac{n}{2} \Big) + \frac{2}{n-1} \sum_{i=1}^{n-1} EX_i. \end{split}$$

Clearly,

$$\mathbb{E}I\left(U_{n-1} \neq \frac{n}{2}\right) = \begin{cases} \frac{n-2}{n-1} & \text{if } n \ge 2, n \mod 2 = 0, \\ 1 & \text{if } n \ge 3, n \mod 2 = 1. \end{cases}$$

Using $x_n = \mathbb{E}X_n$ and $\mathbb{E}I(U_{n-1} \neq \frac{n}{2})$, we may apply Lemma 1 to find

 $\mathbb{E}X_n = a_n$

 $+n \left(\sum_{k=2}^{n-1} \frac{2}{k(k+1)} - \sum_{k=1}^{\left[(n-1)/2 \right]} \frac{2}{(2k-1)2k(2k+1)} \right)$

$$= a_n + n \left(1 - \frac{2}{n} - \sum_{k=1}^{\lfloor (n-1)/2 \rfloor} \frac{1}{2k-1} - \frac{2}{2k} + \frac{1}{2k+1} \right)$$

The first term of the summation in the previous expression is equal to $\frac{1}{2k+1}$ furthermore, for any k 1, we have $\frac{1}{2k+1} = \frac{1}{2(k+1)-1}$ so that the summation can be written as

$$\left(2\left[\frac{n+1}{2}\right]\right)^{-1} + \sum_{k=2}^{\lfloor (n-1)/2 \rfloor} \left(\frac{2}{2k-1} - \frac{2}{2k}\right)$$

$$= \left(2\left[\frac{n+1}{2}\right]\right)^{-1} + 2\sum_{k=3}^{\left[\binom{n-1}{2}\right]} \frac{(-1)^{k+1}}{k}$$

This implies

$$\mathbb{E}[X_n] = a_n + n(1 - \frac{2}{n} - (2\left[\frac{n+1}{2}\right])^{-1} - 2\sum_{k=3}^{(n-1)/2} \frac{(-1)^{k+1}}{k}) = 2n(1 - \ln 2) + O(1) \approx n \quad (9)$$

 $\cdot 0.6137.$

In order to estimate the entropy we recall function Z(t) representing the number of internal vertices of t with isomorphic subtrees. Obviously,

X(t)+Y(t)+Z(t)=n-1

Given $s \in S$ we may define Z (t, s) as the number of internal vertices with both subtrees isomorphic to s. Clearly, $Z(T_n) = \sum_{s \in S} Z(T_n, s)$, and $\mathbb{E}Z_n$: = $\mathbb{E}Z(T_n)$

Let us also use s * s as a shorthand for a non-plane tree having both subtrees of a root isomorphic to s. The stochastic recurrence on $Z_n(s)$ becomes

$$Z_{n}(s) = I(T_{n} \sim s * s) + Z_{U_{n-1}}(s) + Z_{n-U_{n-1}}(s)$$

which leads us to

$$\mathbb{E}Z_n(s) = \mathbb{E}I\left(T_n \sim s * s\right) + \frac{2}{n-1} \sum_{k=1}^{n-1} \mathbb{E}Z_k(s).$$

Moreover, since under the condition that $\Delta(T_n^L) = k$ the event that $T_n \sim s * s$ is equivalent to the intersection of the events $T_n^L \sim s$ and $T_n^R \sim s$, we have

$$\mathbb{E}I(T_n \sim s * s) = \frac{1}{n-1} \sum_{k=1}^{n-1} \mathbb{P}(T_n \sim s * s \mid U_{n-1} = k) = I(n = 2\Delta(s)) \frac{\mathbb{P}^2(T_{n/2} \sim s)}{n-1}$$

Now, we apply Lemma 1 with $x_n = \mathbb{E}Z_n(s)$ and

$$a_n = I(n = 2\Delta(s)) \frac{\mathbb{P}^2(T_{n/2} \sim s)}{n-1}$$

to ultimately find (after some algebra)

$$\mathbb{E}Z_n = n \sum_{k=1}^{[(n+1)/2]} \frac{b_k}{(2k-1)k(2k+1)} + O(1), \quad (10)$$

where $b_k = \sum_{s_k \in \mathcal{T}_k} \mathbb{P}^2(T_k = s_k)$ (see also [21] for analytic derivations of (10)). It is easy to

compute b_k (see Fig. 5) for a few small values of k, namely

$$b_1 = b_2 = 1$$
, $b_3 = \frac{1}{2}$, $b_4 = \frac{2}{9}$, $b_5 = \frac{13}{144}$, $b_6 = \frac{7}{200}$.

In fact, in [21] we show that b_n satisfies for n-2 the following recurrence

$$b_n = \frac{1}{(n-1)^2} \sum_{j=1}^{n-1} b_j b_{n-j}$$

with $b_1 = 1$.

Using (11) we find $\mathbb{E}Z_n \approx n(0.3725 \pm 10^{-4})$, and therefore

$$H(T_n | S_n) = n - 1 - \mathbb{E}[Z_n] \approx n \cdot 0.6275...$$

Since $H(T_n) - H(S_n) = H(T_n|S_n)$, on average the compression of the structure (the non-plane tree) requires asymptotically 0.6275n fewer bits than the compression of any plane tree isomorphic to it. Furthermore, using Theorem 2, we conclude this section with the following second main result.

Theorem 3. The entropy rate $h(s) = \lim_{n \to \infty} H(S_n)/n$ of the non-plane trees is

$$h(s) = h(t) - h(t|s) \approx 1.109...$$

where

$$h(t|s) = 1 - \sum_{k=1}^{\infty} \frac{b_k}{(2k-1)k(2k+1)}$$

$$h(t) = 2\sum_{k=1}^{\infty} \frac{\log_2(k)}{(k+1)(k+2)}$$

with $b_k = \sum_{s_k \in \mathcal{T}_k} \mathbb{P}^2 (T_k = 5_k)$

Remark. The probability $b_k = \sum_{s_k \in \mathcal{T}_k} \mathbb{P}^2(T_k = 5_k)$ is often called the *coincidence probability*,

and in fact is related to the Renyi entropy of order 2 for trees that we introduce next.

Recall that for general a = 0, the Rényi entropy of a random variable X supported on some set X is given by

$$h_{\alpha}(X) = \frac{1}{1-\alpha} \log \left(\sum_{x \in \mathcal{X}} \mathbb{P}[X=x]^{\alpha} \right)$$

We will specialize the above to our distribution on trees and a = 2. It is relatively easy to check that the following quantity is $\Theta(1)$ as $n \to \infty$:

$$h_{2}^{t}(n) = \frac{-\log \sum_{t \in \mathcal{T}_{n}} \mathbb{P}^{2}(T_{n} = t)}{n}$$

and we write $h_2^t = \lim_{n \to \infty} h_2^t(n)$, if the limit exists. Then

$$b_n = \sum_t \mathbb{P}^2 (T_n = t) \operatorname{-exp} (-nh_2^t)$$

for large *n*. Actually, using [22] we prove in [21] that

 $b_n \sim 6n \cdot \rho^n$

where $\rho = 0.3183843834378459 \dots$. Thus $h_2^t = -\log(\rho)$

IV. ALGORITHMS

A. Algorithm for plane trees

An optimal (up to two bits) compression algorithm for plane trees with names can be implemented in two stages: the first compresses the tree T_n , and the second compresses the name function F_n , conditioned on T_n . Both stages can be implemented via arithmetic coding, where each refinement of the interval in the scheme corresponds to a step in a depth-first traversal of the input tree. The result is an algorithm that runs in $O(n^2 \log^2 n \log \log n)$ arithmetic operations in the worst case. As the details for this plane case are quite simple, we leave them to the reader.

Below, we describe in detail the more complicated case of non-plane trees. Again, our schemes rely on arithmetic coding. In this case, we only consider non-plane trees without names to simplify our presentation.

B. Algorithms for non-plane trees

For non-plane trees without vertex names, we present two algorithms: a suboptimal (in terms of expected code length) compression algorithm called COMPRESSNPTREE that runs in worst-case O(n) time, and an optimal algorithm OPTNPTREE that runs in worst-case $O(n^2)$ time (provided we can multiply two binary numbers in O(1) time).

Given the optimal algorithm above, non-plane trees with vertex names may be optimally compressed in a manner similar to plane trees: namely, one first compresses the tree structure, and then compresses the labels conditioned on the structure. We omit the details.

B1: Suboptimal but Time-Efficient Algorithm—We use again arithemtic encoding to compress the tree. We first set the interval to [0,1) and then traverse the vertices of the tree. However, we visit always the smaller subtree (that is, the one with the smaller number of leaves) first (ties are broken arbitrarily).

At each step, if we visit an internal node *v*, we split the interval according to the probabilities of the sizes of the smaller subtree - that is, if the subtree rooted at *v* has *k* leaves and its smaller subtree has *I* leaves, then if *k* is even we split the interval into $\frac{k}{2} - 1$ parts of length $\frac{2}{k-1}$ and one interval of length $\frac{1}{k-1}$. Otherwise, *k* is odd, so we split the interval into parts, all of equal length. For example, if k = 6, then the subintervals have lengths $\frac{2}{5}$, $\frac{2}{5}$ and $\frac{1}{5}$ of the original interval. If k = 7, then the subintervals have length equal to $\frac{2}{6}$ of the original interval. Finally, in both cases we pick *I*-th subinterval as the new interval. The pseudocode of algorithm is presented below:

Function CompressNPTree(*s_n*)

$$[a, b) \leftarrow \text{CompressNPTreeRec}(s_n)$$

$$p \leftarrow b - a, x \leftarrow \frac{a+b}{2}$$

$$C_n^{(2)} = \text{ first } [-\log p] + 1 \text{ bits of } x$$

function CompressNPTreeRec(s)

$$l \leftarrow 0, h \leftarrow 1$$

if $\Delta(s) \ge 2$ then

if
$$\Delta(s^L) \leq \Delta(s^R)$$
 then



else

 $w_1 \leftarrow s^R, w_2 \leftarrow s^L$

$$\left[l_{left}, h_{left}\right] \leftarrow \text{CompressNPTreeRec}(w_1)$$

range $\leftarrow h - l$

 $h \leftarrow l + range * h_{left}$

 $l \leftarrow l + range * l_{left}$

 $\left[l_{\text{right}}, h_{\text{right}}\right] \leftarrow \text{CompressNPTreeRec}(w_2)$

range $\leftarrow h - l$

 $h \leftarrow l + range * h_{right}$

 $l \leftarrow l + range * l_{right}$

return [*1, h*)

Here, we abuse notation slightly: we assume that the non-plane tree s_n is given as an arbitrary plane-oriented representative. Then s_n^L and s_n^R are defined with respect to this representative.

Next, we present a proof of correctness of the COMPRESS- NPTREE algorithm for the nonplane trees, together with the analysis of its performance. The algorithm does not match the entropy rate for non-plane trees, since for every vertex with two non isomorphic subtrees of equal sizes, it can visit either subtree first - resulting in different output intervals, so different

codewords too. Clearly, such intervals are shorter than the length of the optimal interval (which would be equal to the sum of the lengths of all intervals that can be obtained using COMPRESSNPTREE, given a tree s as an input); therefore, the codeword will be longer. However, given EY_n , the expected redundancy rate is within 1% of the entropy rate as proved in Theorem 4 below.

Lemma 2. For a given non-plane tree s with exactly Y(s) vertices with balanced but not isomorphic subtrees, Com- PRESSNPTREE computes an interval whose length is equal to $2^{-Y(s)} \mathbb{P}(S_n = s)$.

Proof: Let $[a_s, b_s)$ be an interval returned by our algorithm, provided its input was *s*. Then, in fact we want to prove that $b_s - a_s = 2^{-Y(s)} \mathbb{P}(S_n = s)$.

Fix a plane tree $t \sim s$ corresponding to the order of the visited vertices of s so that for any vertex of t, its left subtree is visited before its right subtree. Let its interval (compressed using the algorithm for plane trees) be $[a_t, b_t]$.

Now, observe that $b_s - a_s = 2^{X(s)}(b_t - a_t)$, where X (s) is a number of vertices with unbalanced subtrees in *s*, because at each step, when we encounter such vertex, we scale the length of the interval twice as much for *t* when compared to *s*. When we encounter a vertex with balanced subtree, we scale both equally.

However, we already noted that all plane tree isomorphic to *s* have the same probabilities and there are $|[s]| = 2^{X(s) + Y(s)}$ many of them. Therefore (knowing that X(s) = X(t) and Y(s) = Y(t) for any $t \sim s$), we have

$$\mathbb{P}(S_n = s) = |[s]|\mathbb{P}(T_n = t) = 2^{X(s) + Y(s)}\mathbb{P}(T_n = t) = 2^{X(t) + Y(t)}(b_t - a_t) = 2^{Y(s)}(b_s - a_s),$$

which completes the proof.

If we use the arithmetic coding scheme on an interval generated from a tree *s*, we find a codeword $C_n^{(2)}$, which, as a binary number, is guaranteed to be inside this interval. Moreover, we know the following.

Theorem 4. The average length of a codeword $C_n^{(2)}$ from algorithm CompressNPTREE does not exceed 1.013*H*(S_n), where *H*(S_n) is entropy estimated in Theorem 3.

Proof: From Lemma 2, we know that

$$\mathbb{E}C_n^{(2)} = \sum_{s \in \mathcal{S}_n} \mathbb{P}(S_n = s) \left(\left[-\log_2 \left(2^{-Y(s)} \mathbb{P}(S_n = s) \right) \right] + 1 \right)$$

$$< \sum_{s \in S_n} \mathbb{P}(S_n = s) (-\log_2 \mathbb{P}(S_n = s) + Y(s) + 2)$$
$$= H(S_n) + 2 + \sum_{c \in S} \mathbb{P}(S_n = s) Y(s)$$
$$= H(S_n) + 2 + \sum_{t \in \mathcal{T}_n} \mathbb{P}(T_n = s) Y(t) = H(S_n) + 2 + \mathbb{E}Y_n.$$

When combined with the asymptotic behavior of
$$Y_n$$
 we find, using (9) and (10)

$$\lim_{n \to \infty} \frac{\mathbb{E}Y_n}{n} = \lim_{n \to \infty} \frac{n - 1 - \mathbb{E}X_n - \mathbb{E}Z_n}{n} \le 0.014$$

that leads to

$$\lim_{n \to \infty} \frac{\mathbb{E}C_n^{(2)}}{n} \le \lim_{n \to \infty} \frac{H(S_n)}{n} + \frac{\mathbb{E}Y_n}{n} \le 1.124$$

and

$$\lim_{n \to \infty} \frac{\mathbb{E}C_n^{(2)}}{H(S_n)} \le 1 + \lim_{n \to \infty} \frac{\mathbb{E}Y_n}{H(S_n)} \le 1 + \frac{0.014}{1.109} \le 1.013$$

as needed.

B2: Optimal algorithm for non-plane trees—In this section we present an optimal compression algorithm for non-plane trees based on arithmetic encoding. In order to accomplish it we need to define a total order among non-plane trees and compute efficiently the probability distribution $\mathbb{P}(S_n < s)$, where < is the order to be defined. Recall again that *S* is the set of all non-plane trees and S_n is the set of all non-plane trees on *n* leaves. Furthermore, (*s*) is the number of leaves of the non-plane tree *s*.

We start with the definition of our total ordering. In what follows, we will denote by subtrees(s) the set of subtrees of the tree *s* rooted at the children of its root.

Definition 1 (Total ordering on the set of non-plane trees). The relation < on *S* is defined as follows: $s_1 < s_2$ if and only if one of the following holds:

- $(s_1) < (s_2),$
- or $(s_1) = (s_2)$ and min{subtrees (s_1) } < min{subtrees (s_2) },

• or $(s_1) = (s_2)$, min{subtrees (s_1) } = min{subtrees (s_2) } and max{subtrees (s_1) } < max{subtrees (s_2) }.

Here, min and max are defined recursively in terms of the order relation.

Theorem 5. The relation < is a total ordering on *S*.

Proof: The reflexivity and anti-symmetry are straightforward since either $s_1 < s_2$ or $s_1 = s_2$ or $s_1 > s_2$.

To prove the transitivity, we assume that $s_1 < s_2$ and $s_2 < s_3$. Now, if $(s_1) < (s_2)$ or $(s_2) < (s_3)$, then $(s_1) < (s_3)$ (so $s_1 < s_3$), as from the definition of < we know that (s_1) (s_2) (s_3) .

The only remaining possibility is that $(s_1) = (s_2) = (s_3)$. Then, we proceed similarly: if min{subtrees (s_1) } < min{subtrees (s_2) } or min{subtrees (s_2) } < min{subtrees (s_3) }, then min{subtrees (s_1) } < min{subtrees (s_3) } (so $s_1 < s_3$) since from the definition of < if $(s_1) = (s_2) = (s_3)$, then min{subtrees (s_1) } min{subtrees (s_2) } min{subtrees (s_3) }.

Therefore, the only missing case is when $(s_1) = (s_2) = (s_3)$ and min{subtrees (s_1) } = min{subtrees (s_2) } = min{subtrees (s_3) }. Since we know that $s_1 < s_2$ and $s_2 < s_3$, this implies that max{subtrees (s_1) } < max{subtrees (s_2) } < max{subtrees (s_3) } by induction, and this completes the proof.

In what follows, we denote by less(*s*), gtr(*s*) the minimum/maximum root subtree, respectively, of *s* under the ordering just introduced. Moreover, we let T_n denote the (random) plane tree from which S_n is generated, and we recall the notation T_n^L and T_n^R for the left and right subtrees of T_n .

The basic plan is to determine an efficient algorithm that, given a non-plane tree *s*, outputs $\mathbb{P}(S_n < s)$ and $\mathbb{P}(S_n = s)$. This will allow us to construct an arithmetic coding scheme as follows: we associate to s the half-open interval [*a*, *b*), whose left endpoint is given by $\mathbb{P}(S_n < s)$ and whose right endpoint is $\mathbb{P}(S_n \le s)$. The length of this interval is clearly $\mathbb{P}(S_n = s)$, and, because of our total ordering on non-plane trees, for two trees $s_1 < s_2$, the

right endpoint of s_1 is less than or equal to the left endpoint of s_2 . That is, the two intervals do not overlap. Having this interval in hand, we take the midpoint and truncate its binary representation as usual. This will give a uniquely decodable code whose expected length is within 2 bits of the entropy $H(S_n)$.

Now we are in a position to derive the probabilities $\mathbb{P}(S_n = s)$ and $\mathbb{P}(S_n < s)$.

First, we derive an expression for $\mathbb{P}(S_n = s)$ In the case where s has two non-equal subtrees, we have

$$\mathbb{P}(S_n = s) = \mathbb{P}(less(S_n) = less(s), gtr(S_n) = gtr(s)) = \mathbb{P}(T_n^L \sim less(s), T_n^R \sim gtr(s))$$
(12)
+ $\mathbb{P}(T_n^L \sim gtr(s), T_n^R \sim less(s))$

where we recall ~ denotes isomorphism. To calculate the first term on the right-hand side, we condition on the number of leaves in the left subtree of T_n taking the correct value:

$$\mathbb{P}(T_n^L \sim less(s), T_n^R \sim gtr(s))$$

$$= \mathbb{P}\Big(\Delta\Big(T_n^L\Big) = \Delta(less(s))\Big)$$

$$\cdot \mathbb{P}\Big(T_n^L \sim less(s), T_n^R \sim gtr(s)\Big)\Big|\Delta\Big(T_n^L\Big) = \Delta(less(s)))$$

$$\frac{1}{(n-1)} \cdot \mathbb{P}\Big(S_{\Delta}(less(s)) = less(s)\Big) \cdot \mathbb{P}\Big(S_{\Delta}(gtr(s)) = gtr(s))$$

Here, we have applied the conditional independence of the left and right subtrees of T_n given the number of leaves in each. It turns out that the second term of (12) is equal to the first, so we get in this case

=

).

$$\mathbb{P}(\beta_n = s) = \frac{2}{(n-1)} \cdot \mathbb{P}(S_{\Delta(less(s))} = less(s)) \cdot \mathbb{P}(S_{\Delta(gtr(s))} = gtr(s)).$$

In the case where the two subtrees of s are identical, only a single term in (12) is present, and it evaluates to

$$\mathbb{P}(S_n = s) = \frac{1}{n-1} \cdot \mathbb{P}(S_{\Delta(less(s))} = less(s)) \cdot \mathbb{P}(S_{\Delta(gtr(s))} = gtr(s))$$

$$= \frac{1}{n-1} \cdot \mathbb{P}^2 \Big(S_{\Delta(less(s))} = less(s) \Big).$$

Thus, in each case, we have derived a formula for $P(S_n = s)$ that may be recursively computed with O(n) arithmetic operations in the worst case.

It remains to derive an expression for $P(S_n < s)$ We again first consider the case where s has two non-identical subtrees. Following the definition of < , this event is equivalent to

$$[\Delta(less(S_n)) < \Delta(less(s))] \cup [[\Delta(less(S_n)) = \Delta(less(s))] \cap [\Delta(less(S_n)) < \Delta(less(s))]]$$
(13)

$$\cup [[\Delta(less(S_n)) = \Delta(less(s))] \cap [less(S_n) = less(s)] \cap [gtr(S_n) < gtr(s)]]$$

The terms of this union are disjoint, so the total probability is the sum of the probabilities of the individual intersection events.

The first event is equivalent to the union of the disjoint events that the left subtree of T_n has $\langle \Delta(less(s)) |$ leaves or more than n - (less(s)) leaves. The probability of this event is thus

$$\mathbb{P}\left(\Delta\left(less(S_n)\right) < \Delta(less(s))\right) = 2 \cdot \frac{\Delta(less(s)) - 1}{n - 1}$$

The second and third events can be similarly written in terms of disjoint unions of events involving the left and right subtrees of T_n . This gives recursive formulas for their probabilities. Since the formulas are conceptually simple to derive but tedious to write out explicitly, we do not list them, but we mention that at most 4 recursive tree comparison calls are necessary to evaluate $P(S_n < s)$ we need to know the probability that a tree of the appropriate size is < less(s), < gtr(s), and = gtr(s).

Finally, we compute $P(S_n < s)$ in the case where the two subtrees of s are equal. This happens if

$$\Delta\left(less(S_n)\right) < n/2$$

or $\Delta(less(S_n)) = n/2$ and either subtree of S_n is less than the tree s' comprising the two subtrees of s.

The probability of the first event is

$$\begin{split} & \mathbb{P}\left(\Delta\left(less\left(S_{n}\right)\right) < n/2\right) \\ & = \mathbb{P}\left(\Delta\left(T_{n}^{L}\right) < n/2\right) + \mathbb{P}\left(\Delta\left(T_{n}^{L}\right) > n/2\right) \\ & = 1 - \frac{1}{n-1} \end{split}$$

The probability of the second event may be computed by conditioning: the probability that $\Delta(less(S_n)) = \Delta(T_n^L) = n/2$ is $\frac{1}{n-1}$, and the probability, conditioned on this event, that either subtree of S_n is less than s' is

$$\mathbb{P}\Big(T_n^L < s' \cup T_n^R < s' \left| \Delta \Big(T_n^L \Big) = n/2 \Big) = 1 - \mathbb{P}\Big(T_n^L \ge s' \left| \Delta \Big(T_n^L \Big) = n/2 \Big) \cdot \mathbb{P}\Big(T_n^R \ge s' \left| \Delta \Big(T_n^L \Big) = n/2 \Big) = 1 - (1 - \mathbb{P}\big(S_{n/2} < s'\big)\big)^2$$

where we have used the conditional independence of the two subtrees of T_n given their sizes. Thus, the quantities $P(S_n < s)$ and $P(S_n = s)$ may be recursively computed. Importantly, all of the involved probabilities are rational numbers, and arithmetic operations are only performed on integers with value at most $O(2^{n^2})$, so that the interval corresponding to *s* in our arithmetic coding scheme is exactly computable in polynomial time.

With these results in hand, the aforementioned arithmetic coding scheme, in which the interval corresponding to a tree $s \in S_n$ is $[\mathbb{P}(S_n < s), \mathbb{P}(S_n \le s))$, which we call OPT- NPTREE, yields the following optimal compression result:

Theorem 6. The expected code length of the algorithm OPT - NPTREE is at most $H(S_n) + 2$ bits.

Finally, we analyze the running time of the natural recursive algorithm for computing the left and right endpoints of the interval for *s*. We note that, to use the recursive formulas for the probabilities, we need to know less(*s*). To facilitate this, we can first construct a *canonical representative* plane tree *t* isomorphic to *s*, in which for each internal node of *t*, the left subtree is less than or equal to the right one. Note that the left and right subtrees of t are then canonical representations of less(*s*) and gtr(*s*), respectively. To construct such a tree from an arbitrary plane representative of $s \in S_n$ can be done recursively, as the following algorithm shows:

Function CANONIZE (t_n)

> Base case

if n < 2 then return t_n

> Recursive case

CANONIZE (tR)

CANONIZE (tR)

if $t_n^L > t_n^R$ then

swap t^R , t^R

return t_n

When testing whether or not $t_n^L > t_n^R$, we can take advantage of the fact that both subtrees have already been canonized. This saves some time in determining which subtree of each subtree is the lesser one. Nonetheless, the number of comparisons needed to compare two trees is O(n) in the worst case.

To analyze the CANONIZE algorithm, denote by C(n) the number of integer comparisons used by the algorithm on a given tree *t*. It satisfies the recurrence

$$C(n) = C\left(\Delta\left(t^{L}\right)\right) + C\left(\Delta\left(t^{R}\right)\right) + O(n),$$

with a base case of C(1) = O(1). The solution of this recurrence is $O(n^2)$ in the worst case, but on average is $O(n \log n)$.

Having a canonical representative of s in hand, calculating the left and right endpoints of the corresponding interval takes $\Theta(1)$ arithmetic operations at each recursive step (and, thanks to the canonical representation, deciding which subtree is the lesser one takes $\Theta(1)$ time), plus O(n) operations to evaluate the probability $P(S_{\Delta(\text{less}(s))} = \text{less}(s))$ Thus, if we denote by T(n) the number of arithmetic and tree comparison operations to determine the left and right endpoints of the interval for a canonical representation *t* with *n* leaves, we have

$$T(n) \leq T\left(\Delta\left(t^{L}\right)\right) + T\left(\Delta\left(t^{R}\right)\right) + O(n),$$

with a base case of T(1) = O(1). Solving this, we find that, in the worst case, $T(n) = O(n^2)$. Thus, in total, the number of arithmetic and tree comparison operations, including the construction of the canonical representation, is at most $\Theta(n^2)$ in the worst case.

Now we refine the analysis by taking into account the time taken by the arithmetic operations. In the calculation of each interval endpoint, we must keep track of an integer numerator and denominator. In each step, the numerator and denominator are both at most $\Theta(2^{n^2})$, so each may be represented exactly with $\Theta(n^2)$ bits. Each arithmetic operation between two such numbers takes at most $O(n^2 \log n \log \log n)$ bit operations (since multiplying two *N*-digit numbers takes at most $O(N\log N\log \log N)$ operations). Furthermore, taking into account the lengths of the involved integers, the algorithm for computing $P(S_{\Delta(s')} = s')$ takes time $O(n^2 \log^2 n \log \log n)$ (since $\Delta(s') < n$ and an easy inductive proof shows that the lengths of the integers required for the calculation never exceed O(n!)). Then the recurrence for the running time T(n) becomes

$$T(n) = T\left(\Delta\left(t^{L}\right)\right) + T\left(\Delta\left(t^{R}\right)\right) + O\left(n^{2}\log^{2}n\log\log n\right),$$

again with a base case of T(1) = O(1). Solving this, we get that

$$T(n) = O\left(n^3 \log^2 n \log \log n\right)$$

Since this is asymptotically larger than the time taken to construct the canonical representation, the worst-case total running time is $O(n^3 \log^2 n \log \log n)$. The average-case running time can similarly be shown to be $O(n^2 \log^2 n \log \log n)$.

V. CONCLUSION

In this paper we first studied binary plane trees with correlated names and its entropy - which gives the fundamental lower bound on the compression rate - and finally designed an algorithm achieving this bound within an additive constant number of bits in expectation. We also derived the more challenging entropy for non-plane trees and designed an optimal (in terms of expected code length) $O(n^2)$ - time algorithm, as well as a suboptimal O(n)-time algorithm (in a model in which arithmetic operations can be done in O(1) time).

The future directions may focus on construction of universal compression schemes for the presented models (e.g., in the setting where the transition matrix for the names is not known). one may also introduce some horizontal dependency between the letters of the names, using for example mixing sources (noting that a Markov horizontal dependency would be a trivial extension of the results of the present paper). Very recent work [15] concentrates on another nontrivial extension, namely to *d*-ary recursive trees (described in [3]), in which each internal vertex has exactly *d* children.

Finally, as argued in the introduction, our broader goal is to propose adequate models, bounds and algorithms for a wider class of structures, for example random graphs with vertex names.

Acknowledgments

This work was supported by NSF Center for Science of Information (CSoI) Grant CCF-0939370. This work was in addition supported by NSF Grants CCF-1524312, and NIH Grant 1U01CA198941-01.

Biography

Abram Magner is a postdoctoral fellow with the NSF Center for the Science of Information. He conducts research in learning, inference, and data compression problems in network science, random graphs and complex networks, information theory, and random structures, with applications in computer science and biology. He holds BS degrees in mathematics and computer science and MS and Ph.D. degrees in computer science. Prior to graduate school, he has held software engineering intern positions at institutions including IBM and Sandia National Labs. During graduate school, he held a research intern position at the Laboratory for Information, Networking, and Communication Sciences in Paris, France, where he worked on precise average-case analysis of novel tree data structures via methods of complex analysis applied to generating functions.

Krzysztof Turowski is currently Research Scholar at Purdue University. He received his MS and PhD degrees from Gdansk University of Technology, Poland in 2011 and 2015, respectively, both in computer science. From 2010 to 2016 he was employed at the Department of Algorithms and System Modelling at Gdansk University of Technology and from 2016 to 2018 he worked at Google as a software developer for Google Compute Engine. His research interests include graph theory (especially various models of graph coloring), analysis of algorithms and information theory.

Wojciech Szpankowski is Saul Rosen Distinguished Professor of Computer Science at Purdue University where he teaches and conducts research in analysis of algorithms, information theory, analytic combinatorics, data science, random structures, and stability problems of distributed systems. He held several Visiting Professor/Scholar positions, including McGill University, INRIA, France, Stanford, Hewlett-Packard Labs, Universite de Versailles, University of Canterbury, New Zealand, Ecole Polytechnique, France, the Newton Institute, Cambridge, UK, ETH, Zurich, and Gdansk University of Technology, Poland. He is a Fellow of IEEE, and the Erskine Fellow. In 2010 he received the Humboldt Research Award and in 2015 the Inaugural Arden L. Bement Jr. Award. He published two books: "Average Case Analysis of Algorithms on Sequences", John Wiley & Sons, 2001, and 'Analytic Pattern Matching: From DNA to Twitter", Cambridge, 2015. In 2008 he launched the interdisciplinary Institute for Science of Information, and in 2010 he became the Director of the newly established NSF Science and Technology Center for Science of Information.

REFERENCES

- [1]. Ashburner et al., Gene ontology: tool for the unification of biology, Nat. Genet, 2000, 25(1):25–9.[PubMed: 10802651]
- [2]. The Gene Ontology Consortium. Gene ontology consortium: going forward, Nucl. Acids Res, 2015, 43, Database issue D1049–D1056. [PubMed: 25428369]
- [3]. Drmota M, Random Trees: An Interplay between Combinatorics and Probability. Springer Publishing Company, Inc., 2009.
- [4]. Szpankowski W, Average Case Analysis of Algorithms on Sequences. John Wiley & Sons, Inc., New York, NY, 2001.
- [5]. Sedgewick R and Flajolet P, An Introduction to the Analysis of Algorithms. Addison-Wesley Publishing Co., Inc., Reading, MA, 1996.
- [6]. Cover T, Thomas J, Elements of Information Theory, Second Edition, John Wiley & Sons, Inc., Hoboken, NJ, 2005.
- [7]. Naor M, Succinct representation of general unlabeled graphs, Discrete Applied Mathematics, 28(3), 303–307, 1990.
- [8]. Choi Y, Szpankowski W: Compression of Graphical Structures: Fundamental Limits, Algorithms, and Experiments. IEEE Transactions on Information Theory, 2012, 58(2):620–638.
- [9]. Mohri M, Riley M, Suresh AT, Automata and graph compression. ISIT 2015, pp. 2989–2993.
- [10]. Guler B, Yener A, Basu P, Andersen C, Swami A, A study on compressing graphical structures. GlobalSIP 2014, pp. 823–827
- [11]. Turan Gy, On the succinct representation of graphs, Discrete Applied Mathematics, 8(3), 289– 294, 1984.
- [12]. Aldous D and Ross N, Entropy of Some Models of Sparse Random Graphs With Vertex-Names. Probability in the Engineering and Informational Sciences, 2014, 28:145–168.
- [13]. Kieffer JC, Yang E-H, Szpankowski W, Structural complexity of random binary trees. ISIT 2009, pp. 635–639.
- [14]. Zhang J, Yang E-H, Kieffer JC, A Universal Grammar-Based Code for Lossless Compression of Binary Trees. IEEE Transactions on Information Theory, 2014, 60(3):1373–1386.
- [15]. Golebiewski Z, Magner A, Szpankowski W, Entropy of some general plane trees. To appear in Proceedings of the IEEE International Symposium on Information Theory, Aachen, 2017.
- [16]. Steel M, McKenzie A, Distributions of cherries for two models of trees. Mathematical Biosciences, 2000, 164:81–92. [PubMed: 10704639]
- [17]. Steel M, McKenzie A, Properties of phylogenetic trees generated by Yule-type speciation models. Mathematical Biosciences, 2001, 170(1):91–112. [PubMed: 11259805]

- [18]. Abramowitz M, and Stegun I, Handbook of Mathematical Functions, Dover, New York, 1964.
- [19]. Hassani M, Approximation of the dilogarithm function, J. Inequalities in Pure and Applied Mathematics, 8, 1–7, 2007.
- [20]. Bona M, Flajolet P, Isomorphism and symmetries in random phylogenetic trees, Journal of Applied Probability 46(4):1005–1019.
- [21]. Cichon J, Magner A, Szpankowski W, Turowski K, On symmetries of non-plane trees in a nonuniform model. Proceedings of the Fourteenth Workshop on Analytic Algorithmics and Combinatorics, 2017, 156–163.
- [22]. Chern H-H, Fernandez-Camacho M, Hwang H-H, and Martinez C, Psi-series method for equality of random trees and quadratic convolution recurre nces, Random Structures & Algorithms, 44, 67–108, 2012.
- [23]. Peshkin L, Structure induction by lossless graph compression, *In* Proc. of the IEEE Data Compression Conference, 53–62, 2007.
- [24]. Adler M, Mitzenmacher M, Towards Compressing Web Graphs, Data Compression Conference 2001, pp. 203–212.
- [25]. Chierichetti F, Kumar R, Lattanzi S, Mitzenmacher M, Panconesi A, and Raghavan P. On Compressing Social Networks, Proc. ACM KDD, 2009.
- [26]. Matthews SJ, Sul S-J, Williams TL, TreeZip: A New Algorithm for Compressing Large Collections of Evolutionary Trees, Data Compression Conference 2010, pp. 544–554.
- [27]. Sun J, Bollt EM, and Ben-Avraham D, Graph compression-save information by exploiting redundancy, Journal of Statistical Mechanics: Theory and Experiment, P06001, 2008.
- [28]. Schönhage A and Strassen V, Schnelle multiplikation großer zahlen, Computing, 7:281, 1971.



Fig. 1:

Annotated protein interaction network: nodes in the network represent proteins in the human interactome, and there is an edge between two nodes if and only if the two proteins interact in some biological process. Associated with each node is a position in a fixed *ontology* (encoded by a DAG whose nodes represent classes of proteins defined by, say, structural or functional characteristics) known as the Gene Ontology [1], [2]. Nodes that are connected in the interaction network are often close in the ontology.

Author Manuscript



Fig. 2: Linnaean taxonomy of *hominoid*







Fig. 4:

A rooted plane tree and its standardization. The left tree can be represented by the list of triples $\{(5,1,7), (1, 2,3)\}$. After standardization, this becomes $\{(1,2,5), (2,3,4)\}$. Note that this is distinct from the tree $\{(1, 2, 3), (3, 4, 5)\}$ or, equivalently, $\{(5,1, 7), (7, 2, 3)\}$.



