

# The Height of List-tries and TST

N. Broutin<sup>1</sup> and L. Devroye<sup>1†</sup>

<sup>1</sup>*School of Computer Science, McGill University, 3480 University Street, H3A 2K6 Montreal, Canada.*  
*Email: {nbrout, luc}@cs.mcgill.ca.*

*received 14 Oct 1998, revised 15<sup>th</sup> February 2007, accepted tomorrow.*

---

We characterize the asymptotics of heights of the trees of de la Briandais (1959) and the ternary search trees (TST) of Bentley and Sedgewick (1997). Our proof is based on a new approach of the structure of tries that distinguishes the bulk of the tree, called the *core*, and the long trees hanging down the core, called the *spaghettis*.

**Keywords:** Tries, branching process, height, de la Briandais, TST.

---

## 1 Introduction

Tries are data structures used to manipulate and store strings by taking advantage of the digital character of words. They were introduced by de la Briandais (1959). Apparently, the term of *trie* was coined by Fredkin (1960) as a part of the word “retrieval”. Their properties and uses are reviewed by Knuth (1973) and more recently by Szpankowski (2001). Consider  $n$  sequences of characters (or strings) from an alphabet  $\mathcal{A} = \{1, \dots, d\}$ . Each one of the sequences carves a path in an infinite rooted  $d$ -ary position tree  $T_\infty$  where the children of each node are labeled with the characters of  $\mathcal{A}$ : starting from the root, the characters are used in order to move down the tree. If all the sequences are distinct, the corresponding paths in  $T_\infty$  are distinct as well. The trie  $T_n$  is defined to be the smallest subtree of  $T_\infty$  such that the paths corresponding to the sequences are distinct within  $T_n$ .

In this paper, we are interested in tries built from  $n$  independent sequences. Each sequence is an infinite sequence of independent and identically distributed (i.i.d.) characters distributed like  $A$ , where  $\mathbf{P}\{A = i\} = p_i$ . We assume without loss of generality that  $1 > p_1 \geq p_2 \geq \dots \geq p_d > 0$ . The quantity of interest under this model is

$$Q = \sum_{i=1}^d p_i^2, \tag{1}$$

It is well known that the height  $H_n$  of a trie built from  $n$  independent such sequences satisfies (Régnier, 1981; Devroye, 1984; Pittel, 1985; Szpankowski, 1991, 2001)

$$\frac{H_n}{\log n} \xrightarrow{n \rightarrow \infty} \frac{2}{\log(1/Q)} \quad \text{in probability.} \tag{2}$$

---

<sup>†</sup>Research of the authors was supported by NSERC Grant A3456 and a James McGill fellowship.

The trie is only an *abstract data structure*, that is, it does not specify the implementation (see Clément et al., 1998, 2001). The usual implementation of a trie uses an array for the branching structure of a node (Fredkin, 1960). Although this always ensures constant-time shunting of the strings, the space required may become an issue for large alphabets: many pointers would be left unused. To avoid this, one can replace the array by variable size structures. de la Briandais (1959) proposed to use linked-lists, and we shall call the implementation a *list-trie*. More recently, Bentley and Sedgewick (1997) proposed an elegant structure based on binary search trees. It is known as the *bst-trie*, ternary search trie or TST for short.

These structures aim at a trade-off between the storage space and the speed, and the access time to children is no longer constant. In particular, the height of the tree and the worst-case search time are different in general. List-tries and the TST may be seen as *high-level* tries whose edges are weighted to reflect the internal *low-level* structure of a node (see Figure 2). This point of view has been taken by Clément, Flajolet, and Vallée (1998, 2001) who analyzed thoroughly these *hybrid* implementations of tries. In particular, they analyzed the average size and average depth. The question of the worst-case search time in hybrid-tries was left open. The worst-case search time is then given by the *weighted height* of the tree, we shall call it the height in following. This paper addresses the question of the (weighted) height of hybrid tries, and of list-tries and TST, in particular.

## 2 Weighted tries

### 2.1 An embedding to construct weighted tries

In this section we propose an embedding to construct the weighted tries. Strictly speaking, since we are interested in the weighted height, we only construct the sequence of weighted depths of the nodes. Note that our embedding is only *one* way to build tries with the desired distribution. We will see in section 4 that hybrid tries, like list-tries and TST, can be seen as weighted tries. Our construction emphasizes an underlying structure consisting of independent random variables. However, in the coupled tries built from the embedding, the random variables are dependent in general. Consider the distribution  $\{p_1, \dots, p_d\}$  over the alphabet  $\mathcal{A}$ . We assume without loss of generality that  $1 > p_1 \geq p_2 \geq \dots \geq p_d > 0$ . We are given  $n$  independent strings, each consisting of an infinite sequence of i.i.d. characters of  $\mathcal{A}$  distributed as  $A$ , where  $\mathbf{P}\{A = i\} = p_i$ ,  $1 \leq i \leq d$ .

THE SHAPE OF THE TRIE. The shape of the trie is simply an ordinary trie: Each string defines an infinite path in  $T_\infty$ . The *cardinality*  $N_u$  of a node  $u \in T_\infty$  is the number of strings whose path in  $T_\infty$  intersect  $u$ . Then, the trie  $T_n$  is constructed by pruning  $T_\infty$  down to every node of cardinality at most one. The sequences are distinct with probability one, and the strings define distinct paths in  $T_\infty$ . Therefore, the trie  $T_n$  is almost surely finite. The tree  $T_n$  constitutes the *shape* of the weighted trie. We define  $E_i = -\log p_i$ . For the edge  $e$  between  $u$  and its  $i$ -th child in  $T_\infty$ , we let  $p_e = p_i$  and  $E_e = -\log p_e$ .

DIFFERENT TYPES OF NODES. There are  $2^d$  types of nodes, each type being characteristic of the *branching structure* of the node. The branching structure of every node  $u \in T_n$  is described by a  $d$ -vector  $\tau_u$ : if  $u_1, \dots, u_d$  are the  $d$  children of  $u$ , then we define

$$\tau_u = (\mathbf{1}[N_{u_1} \geq 1], \mathbf{1}[N_{u_2} \geq 1], \dots, \mathbf{1}[N_{u_d} \geq 1]).$$

The vector  $\tau_u$  indicates which one of the  $d$  edges down  $u$  are part of some path in  $T_n$ , and influences the cost of accessing the children of  $u$ .

THE WEIGHTS. Consider a sequence of random vectors  $\{Z^\tau : \tau \in \{0, 1\}^d\}$ , where  $Z^\tau = (Z_1^\tau, \dots, Z_d^\tau)$ . The vectors  $Z^\tau$  describe the cost of accessing the children of the nodes. We assume that

- for all  $\tau \in \{0, 1\}^d$ ,  $Z^\tau \leq d$ , and
- for all types  $\tau$  such that  $\tau$  is a permutation of  $(1, 0, \dots, 0)$ ,  $Z^\tau = (1, \dots, 1)$ .

Each node of  $T_\infty$  is assigned an independent copy of the whole sequence. Consider a node  $u \in T_\infty$ , and its sequence  $\{Z^\tau\}$ . Weights are associated with the edges to  $u$ 's children based on its type  $\tau_u$ . The edge  $e_i$  between  $u$  and its  $i$ -th child in  $T_\infty$  is given the weight

$$Z_{e_i} = Z_i^{\tau_u} = \sum_{\tau \in \{0,1\}^d} Z_i^\tau \cdot \mathbf{1}[\tau_u = \tau].$$

We use the notations  $Z_i^\tau$  and  $Z_e$  interchangeably. It should always be clear whether a subscript refers to an index or an edge. Let  $\pi(u)$  be the set of edges on the path from  $u$  up to the root in  $T_\infty$ . The *weighted depth* of a node  $u$  is defined by  $D_u = \sum_{e \in \pi(u)} Z_e$ . A tree that may be constructed by this procedure is called a *weighted trie*. We are interested in the weighted height of  $T_n$ :

$$H_n = \max\{D_u : u \in T_n\} = \max\{D_u : N_u \geq 2\} + 1.$$

Surprisingly, the first asymptotic term of  $H_n$  depends on two parameters only: the distribution  $\{p_1, \dots, p_d\}$ , and  $Z = Z^{(1, \dots, 1)}$ . In particular, the first order asymptotics of  $H_n$  stays the same if we modify  $\{Z^\tau, \tau \in \{0, 1\}^d\}$  in such a way that  $Z$  remains unchanged. This may be explained using the structure of a trie.

## 2.2 The structure of a trie

In a related project, we have explained the profile of a trie by distinguishing a so-called *core*, that constitutes the bulk of the trie, and *spaghetti*-like trees hanging down the core (Broutin and Devroye, 2007a).

THE CORE OF A TRIE. What we call the core here should not be confused with the graph-theoretic core, which happens to be empty for trees (see e.g., Janson et al., 2000). The core of the trie is defined to be the set of nodes  $u \in T_\infty$  for which  $N_u \geq m(n)$ , for  $m(n) \rightarrow \infty$  and  $m(n) = o(\log n)$ . The core is denoted by  $C$ . Since  $m \rightarrow \infty$ , the nodes in the core are of the same type  $\tau = (1, \dots, 1)$  with probability  $1 - o(1)$ . As a consequence, in a weighted trie, the distribution of weights in the core should be closely approximated by  $Z = Z^{(1, \dots, 1)}$ . The core can be described by its *logarithmic profile*

$$\phi(\alpha, t) = \lim_{n \rightarrow \infty} \frac{\log \mathbf{E}P_m(t \log n, \alpha \log n)}{\log n} \quad \forall t, \alpha > 0, \quad (3)$$

where  $P_m(k, h)$  denotes the number of nodes  $u$ ,  $k$  levels away from the root with  $N_u \geq m(n)$  and  $D_u \geq h$ . In other words, assuming for now that the limit in (3) exists, we have  $\mathbf{E}P_m(t \log n, \alpha \log n) = n^{\phi(\alpha, t) + o(1)}$ , as  $n \rightarrow \infty$ . This will be proved, and the function  $\phi(\cdot, \cdot)$  will be characterized in Theorem 2.

HANGING SPAGHETTIS. The *spaghettis* are the trees remaining when pulling out the core from the trie. They lie in the part of the trie where the nodes do not have  $d$  children any more: the types of the nodes take all the values in  $\{0, 1\}^d$  with positive probability. However, the weighted height of a *long* spaghetti is close to the weighted height its unweighted height, or number of levels. To see this, observe that the nodes

not in the core have cardinality at most  $m(n) = o(\log n)$ , so each spaghetti stores at most  $m(n)$  sequences. Each time the type  $\tau$  is not a permutation of  $(1, 0, \dots, 0)$ , i.e., the node is truly branching, at least one string is put aside from the longest path. This can happen at most  $o(\log n)$  times, and hence the heights with and without the branching nodes differ by at most  $o(\log n)$ . If the number of levels is  $\Theta(\log n)$ , as is the case for the highest ones, the difference is negligible.

Both the core and the spaghetti contribute significantly to the height of a weighted trie. By figuring out what the core looks like, we can determine *when* the spaghetti take over. Roughly speaking, we then know if an edge's weight can be approximated by a component of  $\mathcal{Z}$  or simply remain unweighted. The main result of this paper is the following theorem.

**Theorem 1.** *Consider an weighted trie built from  $n$  independent sequences consisting of i.i.d. characters distributed as  $A$ , where  $\mathbf{P}\{A = i\} = p_i$ ,  $1 \leq i \leq d$ . Let  $H_n$  be its weighted height. Let  $\phi(\alpha, t)$  be the logarithmic weighted profile of the core of  $T_n$ . Let*

$$c = \sup \left\{ \alpha + \frac{\phi(\alpha, t)}{\log(1/Q)} : \phi(\alpha, t) \geq 0 \right\},$$

where  $Q = \sum_{i=1}^d p_i^2$ . Then  $H_n = c \log n + o(\log n)$  in probability, as  $n \rightarrow \infty$ .

**Remark.** The respective contributions of the core and spaghetti are  $\alpha$  and  $\phi(\alpha, t)/\log(1/Q)$ .

### 3 The height of weighted tries

#### 3.1 The shape of the core

Consider a weighted trie defined as in section 2. We consider  $m = m(n) \rightarrow \infty$  with  $m(n) = o(\log n)$ . Let  $\mathcal{L}_k$  be the set of nodes  $k$  levels away from the root in  $T_\infty$ . Let  $P_m(k, h)$  be the number of nodes  $u \in \mathcal{L}_k$  with  $D_u \geq h$  and  $N_u \geq m$ . Since  $m \rightarrow \infty$ , for  $n$  large enough, we have  $m \geq b$  and

$$P_m(k, h) = \sum_{u \in \mathcal{L}_k} \mathbf{1}[N_u \geq m, D_u \geq h].$$

The asymptotic properties of the expected profile are directly tied to large deviation theory (Dembo and Zeitouni, 1998). The random vector of interest here is  $(Z, E) = (Z_K, -\log p_K)$ , where  $K$  is uniform in  $\{1, \dots, d\}$ . For  $\lambda, \mu \in \mathbb{R}$ , the associated *generating function of the cumulants* is

$$\Lambda(\lambda, \mu) = \log \mathbf{E} \left[ e^{\lambda Z + \mu E} \right].$$

Then, we define the convex dual  $\Lambda^*(\cdot, \cdot)$  of  $\Lambda$  by, for  $x, y \in \mathbb{R}$ ,

$$\Lambda^*(x, y) = \sup_{\lambda, \mu} \{ \lambda x + \mu y - \Lambda(\lambda, \mu) \}.$$

**Theorem 2.** *Let  $m = m(n) \rightarrow \infty$  with  $m = o(\log n)$ . Let  $k \sim t \log n$  and  $h \sim \alpha \log n$  for some positive constants  $t$  and  $\alpha$ . Let*

$$\phi(\alpha, t) = t \log d - t \cdot \inf \left\{ \Lambda^*(x, y) : x > \frac{\alpha}{t}, y < \frac{1}{t} \right\}. \quad (4)$$

If  $\phi(\alpha, t) > 0$ , then  $\mathbf{E}P_m(k, h) = n^{\phi(\alpha, t) + o(1)}$ , as  $n \rightarrow \infty$ .

**Remarks.** (a) Observe that Theorem 2 justifies the definition of  $\phi(\cdot, \cdot)$  in (3).

(b) The constraint that  $m(n)$  is  $o(\log n)$  is only used in the lower bound. Intuitively, we do not want to shave off too many layers of the tree for the lower bound. This has no effect on the upper bound.

Unlike the profile of *unweighted* tries (Devroye, 2002, 2005; Park et al., 2006), that of *weighted* tries does not seem concentrated. However, it is log-concentrated in the sense of the following theorem, and the true profile  $P_m(k, h)$  is close enough to  $\mathbf{E}P_m(k, h)$ :

**Theorem 3.** *Let  $m = m(n) \rightarrow \infty$  as  $n \rightarrow \infty$  such that  $m = o(\log n)$ . Let  $k \sim t \log n$  and  $h \sim \alpha \log n$  for some positive constants  $t$  and  $\alpha$ . Then, for all  $\varepsilon > 0$ , as  $n \rightarrow \infty$ ,*

$$\mathbf{P} \left\{ P_m(k, h) \leq n^{\phi(t, \alpha) - \varepsilon} \right\} \xrightarrow{n \rightarrow \infty} 0, \quad \text{and} \quad \mathbf{P} \left\{ P_m(k, h) \geq n^{\phi(t, \alpha) + \varepsilon} \right\} \leq n^{-\varepsilon + o(1)}.$$

### 3.2 How long is a spaghetti?

The behavior of the spaghettis is radically different from the one observed in the core. This is because the number of strings stored by a single spaghetti is at most  $m(n)$ . We first look at the profile, not of a single trie, but of a forest of independent tries.

Let  $T^1, T^2, \dots, T^n$  be  $n$  independent tries. We assume that  $T^i$  is a weighted trie on  $m_i = m_i(n)$  sequences generated by a memoryless source with distribution  $\{p_1, \dots, p_d\}$ . Also, we assume that for all  $i$ ,  $m/d \leq m_i \leq m$ . The roots of  $T^i$ ,  $1 \leq i \leq n$ , all lie at level zero. Then, we let  $P^s(k, h)$  count the number of nodes  $u$  at level  $k$  with  $D_u \geq h$  lying in any  $T^i$ . Since  $T^i$  is a trie, we only count the nodes for which  $N_u \geq 2$ . We are interested in  $\mathbf{E}P^s(k, h)$  when  $k \sim \rho \log n$  and  $h \sim \gamma \log n$ . Recall that the difference between the number of layers and the weighted height is at most  $dm(n) = o(\log n)$ . Hence,  $\mathbf{E}P_m(k, h) = o(1)$  if  $h \geq k + \Omega(\log n)$ . When  $h \leq k + o(\log n)$ , only the number of level is relevant. For  $n$  large enough that  $m/d \geq 2$ , by linearity of expectation, we have

$$nQ^k \leq \mathbf{E}P_m(k, h) \leq nm^2Q^k,$$

since two strings are identical up to the  $k$ -th character with probability  $Q^k$ . In other words, we have

**Theorem 4.** *Let  $T^i$ ,  $1 \leq i \leq n$ , be a forest of  $n$  independent tries. Let  $T^i$  store  $m_i = m_i(n)$  sequences. Assume that  $m/d \leq m_i \leq m$  for all  $1 \leq i \leq n$ . Let  $k \sim \rho \log n$  and  $h \sim \gamma \log n$ , as  $n \rightarrow \infty$ , for positive constants  $\rho$  and  $\gamma$ . Then,*

$$\mathbf{E}P^s(k, h) = \begin{cases} n^{1 + \rho \log Q + o(1)} & \text{if } \rho \geq \gamma \\ o(1) & \text{otherwise.} \end{cases}$$

as  $n \rightarrow \infty$ .

So the logarithmic profile of our forest is  $1 + \rho \log Q$ . Theorem 5 claims that the deepest node to occurs when this logarithmic profile vanishes. Let  $H^1, \dots, H^n$  be the weighted heights of  $T^1, \dots, T^n$ , respectively, and define  $S_n = \max\{H^i : 1 \leq i \leq n\}$ . Then:

**Theorem 5.** *Assume that  $p_1 < 1$ . Assume that  $m(n) \rightarrow \infty$  and  $m(n) = o(\log n)$ . Let Then,  $S_n \sim \log_{1/Q} n$  in probability, as  $n \rightarrow \infty$ . Furthermore, for every  $\varepsilon > 0$ , there exists  $\delta > 0$  such that, as  $n \rightarrow \infty$ ,*

$$\mathbf{P} \left\{ \frac{S_n}{\log n} \geq \frac{1}{\log(1/Q)} + \varepsilon \right\} = O(n^{-\delta}). \quad (5)$$

### 3.3 Projecting the profile of the core

In this section, we give a sketch of the proof of Theorem 1. Consider a weighted trie  $T_n$ . The spaghetti are rooted at a node  $u \in \partial C$ , the external node-boundary of the core  $C$  in  $T_n$  (the nodes  $u \in \partial C$  are the children of some node  $v$  in the core, but are not themselves in the core). Then, if we write  $W_u$  for the height of the subtree rooted at  $u$ , we have

$$H_n = \max\{D_u + W_u : u \in \partial C\} = \max_{h,k}\{h + W_u : u \in \partial C \cap \mathcal{L}_k, D_u \geq h\},$$

where the nodes in  $\partial C$  have been split into groups depending on their level  $k$  and weighted depth  $h$ . Then, we can rewrite

$$H_n = \max_{h,k}\{h + \max\{W_u : u \in \partial C, D_u \geq h, u \in \mathcal{L}_k\}\}.$$

We have thus separated the contributions of the core from that of the spaghetti,

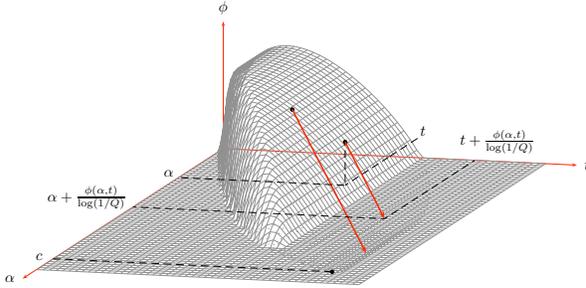
$$h \quad \text{and} \quad \max\{W_u : u \in \partial C, D_u \geq h, u \in \mathcal{L}_k\},$$

respectively. Only the latter expression requires more thought. The set of nodes  $u \in \partial C \cap \mathcal{L}_k$  with  $D_u \geq h$  roughly contains roughly  $n^{\phi(\alpha,t)+o(1)}$  by Theorem 2. Hence  $\max\{W_u : u \in \partial C, D_u \geq h, u \in \mathcal{L}_k\}$  is  $S_{n'}$  with  $n' = n^{\phi(\alpha,t)+o(1)}$ , that is,  $\phi(\alpha,t)/\log(1/Q)$  by Theorem 5. This explains why  $H_n \sim c \log n$ , where

$$c = \sup \left\{ \alpha + \frac{\phi(\alpha,t)}{\log(1/Q)} : \phi(\alpha,t) \geq 0 \right\}.$$

Complete proofs of the theorems can be found in Broutin (2007) or Broutin and Devroye (2007b).

Theorem 1 can be interpreted by projecting the logarithmic profile  $\phi(\alpha,t)$  on the horizontal plane going through the origin. See Figure 1.



**Fig. 1:** A geometric interpretation for the height: each point  $(\alpha, t, \phi(\alpha, t))$  of the logarithmic profile of the core throws a line whose direction is given by  $(1, 1, \log Q)$ . The line intersects the plane  $\phi = 0$  at  $(\alpha - \phi(\alpha, t)/\log Q, t - \phi(\alpha, t)/\log Q, 0)$ . The constant  $c$  is the largest coordinate of one of these point measured along the  $\alpha$ -axis.

## 4 Application: the height of hybrid tries

Let  $\mathcal{A} = \{1, \dots, d\}$  be the alphabet. Let  $\{A^i, 1 \leq i \leq n\}$  be the  $n$  strings. In the following, we distinguish the *nodes* that constitute the high-level trie structure from the *slots* which make the low-level structure of a node, whether this latter be a linked-list or a binary search tree.

The high-level tree between the nodes of an hybrid trie is just the ordinary trie. Only the low-level structure of the nodes differ from the array-based implementation. Consider a node  $u \in T_\infty$ . The subtree

rooted at  $u$  stores a subset of the strings  $A^i$ ,  $1 \leq i \leq n$ . Let  $\mathcal{N}_u \subset \{1, \dots, n\}$  be the set of their indices, and write  $N_u = |\mathcal{N}_u|$ . As in ordinary tries, for a node  $u$  at level  $k$  in  $T_\infty$ , only the  $k$ -th characters of each sequence are used. Moreover, only set  $\mathcal{A}_u \subset \mathcal{A}$  matters, which consists of the characters appearing at the  $k$ -th position in the sequences  $A^i$ ,  $i \in \mathcal{N}_u$ . In an hybrid trie, the order in which the strings are used to build the trie matters, and we let  $\sigma_u$  be the permutation  $\mathcal{A}_u$  where the characters are ordered by first appearance. The internal structure of the node  $u$  is built by successive insertions of the elements of  $\sigma_u$  into an originally empty linked list, or binary search tree. For list-tries we then have a random linked list on  $\mathcal{A}_u$  and for TST a binary search tree on  $\mathcal{A}_u$ . This is shown in Figure 2. The costs of branching to a subtree is then given by the number of edges to cross *inside the node* to the subtree. The distribution clearly depends on the type  $\tau_u$  of node  $u$ , and when  $|\mathcal{A}_u| = 1$ , the cost is simply 1. This explains why hybrid tries are weighted tries in the sense of section 2.



**Fig. 2:** The different node structures used for the standard (top-left), list (bottom-left) and bst-trie (right) when the order of appearance of the characters is 3, 5, 4, 1, 2 and 6. The dashed arrows represent the pointers to further levels of the trie.

### 4.1 list-tries

In the list-trie of de la Briandais (1959), the cost of branching to a character  $a$  is just the index of  $a$  in the permutation  $\sigma_u$ . For every node  $u$ , for which  $\mathcal{A}_u = \mathcal{A}$ ,  $\sigma_u$  is distributed as the sequence (in order) of first appearance of characters in an infinite string generated by the source. This fully describes the distribution of  $Z = Z^{(1, \dots, 1)}$ . Then  $Z_i$  is the index of  $i$  in  $\sigma$ , and  $(Z, E) = (Z_K, -\log p_K)$ , where  $K$  is uniform in  $\{1, \dots, d\}$ .

**Theorem 6.** *Let  $H_n$  be the weighted height of a list-trie on  $n$  sequences generated by the memoryless source with probabilities  $\{p_1, \dots, p_d\}$ . Then,  $H_n \sim c \log n$  in probability, as  $n \rightarrow \infty$ , where*

$$c = \sup_{\alpha > 0} \left\{ \alpha + \frac{\phi(t, \alpha)}{\log(1/Q)} : \phi(\alpha, t) \geq 0 \right\},$$

and  $\phi(\cdot, \cdot)$  is the logarithmic profile of the trie weighted with  $(Z, E)$  is described above.

It seems difficult to obtain  $\phi(\alpha, t)$  for general  $\{p_1, \dots, p_d\}$ . We treat completely a concrete example.

**Example: symmetric list-tries.** For any  $\lambda, \mu \in \mathbb{R}$ , we have

$$\Lambda(\lambda, \mu) = \log \mathbf{E} \left[ e^{\lambda Z} \cdot d^\mu \right] + \mu \log d = \log \left( \sum_{i=1}^d e^{i\lambda} \right) + (\mu - 1) \log d.$$

Then,  $\Lambda^*(x, y) = \sup_{\lambda, \mu} \{\lambda x + \mu y - \Lambda(\lambda, \mu)\}$ . For  $x \in [1, d]$ , there exists  $\lambda = \lambda(x)$  such that

$$x = \frac{\partial \Lambda(\lambda, \mu)}{\partial \lambda} = \frac{\sum_{i=1}^d i e^{i\lambda}}{\sum_{i=1}^d e^{i\lambda}}. \quad (6)$$

Therefore, we have

$$\Lambda^*(x, y) = \begin{cases} \lambda x - \log(\sum_{i=1}^d e^{i\lambda}) + \log d & \text{if } x \in [1, d], y = \log d \\ \infty & \text{otherwise.} \end{cases}$$

For instance, for  $d = 2$ , we can find an explicit expression for  $\Lambda^*$ : for  $x \in [1, d]$ ,

$$\Lambda^*(x, \log 2) = x \log \left( \frac{1-x}{x-2} \right) - \log \left( \frac{1-x}{2-x} + \left( \frac{1-x}{2-x} \right)^2 \right) + \log 2.$$

Numerical values for  $c = c(d)$  can be found in Table 1.

d	2	3	10
c	3.28661...	3.12515...	4.92852...

**Tab. 1:** Some numerical values of  $c = c(d)$  characterizing the height of symmetric list-tries.

## 4.2 Ternary search trees

In the ternary search trees introduced by Bentley and Sedgewick (1997), the implementation of a node uses a binary search tree. The cost of branching to a character  $i \in \mathcal{A}$  at a node  $u$  is one plus the depth of  $i$  in the binary search built from the (non-uniform) random permutation  $\sigma_u$ . When the node  $u$  is of type  $\tau_u = (1, \dots, 1)$ , the permutation  $\sigma_u$  is distributed as  $\sigma$ , the ordered list of first appearances of characters in an infinite string generated by the memoryless source with distribution  $\{p_1, \dots, p_d\}$ .

Let  $Z_i$  be distributed as the depth of  $i$  in the binary search tree built from  $\sigma$ . Then,  $Z$  is distributed as  $(Z_1, \dots, Z_d)$  and  $(Z, E) = (Z_K, -\log p_K)$ , where  $K$  is uniform in  $\{1, \dots, d\}$ . By Theorem 1, we obtain:

**Theorem 7.** *Let  $H_n$  be the weighted height of a TST on  $n$  sequences. Let*

$$c = \sup \left\{ \alpha + \frac{\phi(\alpha, t)}{\log(1/Q)} : \phi(\alpha, t) \geq 0 \right\},$$

where  $\phi(\alpha, t)$  is the logarithmic profile of the core of a trie weighted by  $(Z, E)$  described above. Then,  $H_n \sim c \log n$  in probability, as  $n \rightarrow \infty$ .

The random vector  $(Z, E)$  seems complicated to describe for general distributions  $p_1, \dots, p_d$ . Some parameters like the average value and the variance of  $Z_i$ ,  $1 \leq i \leq d$ , have been studied by Clément et al. (1998, 2001) and Archibald and Clément (2006).

**Example: Symmetric TST.** We assume here that  $p_1 = p_2 = \dots = p_d$ . In this case, the permutation  $\sigma$  is just a uniform random permutation. Hence,  $Z_i$  is the depth of the key  $i$  in a random binary search tree.

Observe that unlike in the case of list-tries,  $Z_i$ ,  $1 \leq i \leq d$ , do *not* have the same distribution. This is easily seen, since, for instance as  $d \rightarrow \infty$ ,  $\mathbf{E}Z_1 \sim \log d$  whereas  $\mathbf{E}Z_{\lfloor d/2 \rfloor} \sim 2 \log d$ . However, we are only interested in the distribution of  $Z$ , that is, the depth of a uniform random node. This distribution is known exactly, and is due to Brown and Shubert (1984):

$$\mathbf{P}\{Z = k\} = \frac{2^{d-1}}{d \cdot d!} \sum_{j=k}^d \begin{bmatrix} d \\ j \end{bmatrix}, \quad (7)$$

where  $\begin{bmatrix} n \\ k \end{bmatrix}$  denotes the Stirling number of the first kind with parameter  $n$  and  $k$  (see Sedgewick and Flajolet, 1996; Mahmoud, 1992). Using (7), it is possible to compute the cumulant generating function  $\Lambda$ , and  $\phi(\alpha, t)$ . Observe that when  $d = 2$ , TST are equivalent to list-tries. When  $d = 3$ , we have

$$Z = \begin{cases} 1 & \text{w.p. } 1/3 \\ 2 & \text{w.p. } 1/2 \\ 3 & \text{w.p. } 1/6. \end{cases}$$

Numerical values for the constant  $c = c(d)$  such that  $H_n \sim c \log n$  in probability as  $n \rightarrow \infty$  are given in Table 2.

d	2	3	10
c	3.28661...	2.90777...	

**Tab. 2:** Some numerical values of  $c = c(d)$  characterizing the height of symmetric ternary search trees.

## 5 Concluding remarks

Theorem 1 is not the most general one can obtain. In particular, the constraints that the components of  $Z^\tau$  are bounded may be replaced by bounds on the moment generating function. We weights for non-branching nodes may also be random instead of the unit values used here. Finally, similar results hold for  $b$ -tries, where the leaves are allowed to contain up to  $b$  strings. These extensions and complete proofs of the theorems can be found in Broutin (2007) or Broutin and Devroye (2007b).

## 6 Acknowledgement

We are very grateful to Julien Clément for bringing our attention on this problem and for his insightful comments.

## References

- M. Archibald and J. Clément. Average depth in binary search tree with repeated keys. In *Fourth Colloquium on Mathematics and Computer Science*, 2006.
- J. L. Bentley and R. Sedgewick. Fast algorithm for sorting and searching strings. In *Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 360–369, 1997.

- N. Broutin. *Shedding new light on random trees*. PhD thesis, McGill University, Montreal, QC, 2007.
- N. Broutin and L. Devroye. The core of a trie. Manuscript, February 2007a.
- N. Broutin and L. Devroye. Weighted height of random tries. Manuscript, February 2007b.
- G.G. Brown and B.O. Shubert. On random binary trees. *Mathematics of Operations Research*, 9:43–65, 1984.
- J. Clément, P. Flajolet, and B. Vallée. The analysis of hybrid trie structures. In *9th annual ACM-SIAM Symposium on Discrete Algorithms*, pages 531–539, Philadelphia, PA, 1998. SIAM Press.
- J. Clément, P. Flajolet, and B. Vallée. Dynamical source in information theory: a general analysis of trie structures. *Algorithmica*, 29:307–369, 2001.
- R. de la Briandais. File searching using variable length keys. In *Proceedings of the Western Joint Computer Conference, Montvale, NJ, USA*. AFIPS Press, 1959.
- A. Dembo and O. Zeitouni. *Large Deviation Techniques and Applications*. Springer Verlag, second edition, 1998.
- L. Devroye. Laws of large numbers and tail inequalities for random tries and PATRICIA trees. *Journal of Computational and Applied Mathematics*, 142:27–37, 2002.
- L. Devroye. Universal asymptotics for random tries and PATRICIA trees. *Algorithmica*, 42:11–29, 2005.
- L. Devroye. A probabilistic analysis of the height of tries and of the complexity of triesort. *Acta Informatica*, 21: 229–237, 1984.
- E. Fredkin. Trie memory. *Communications of the ACM*, 3(9):490–499, 1960.
- S. Janson, T. Łuczak, and A. Ruciński. *Random Graphs*. Wiley, New York, 2000.
- D. E. Knuth. *The Art of Computer Programming: Sorting and Searching*, volume 3. Addison-Wesley, Reading, MA, 1973.
- H. Mahmoud. *Evolution of Random Search Trees*. Wiley, New York, 1992.
- G. Park, H.K. Hwang, P. Nicodème, and W. Szpankowski. Profile of tries. manuscript, 2006.
- B. Pittel. Asymptotic growth of a class of random trees. *The Annals of Probability*, 13:414–427, 1985.
- M. Régnier. On the average height of trees in digital search and dynamic hashing. *Information Processing Letters*, 13:64–66, 1981.
- R. Sedgewick and P. Flajolet. *An Introduction to the Analysis of Algorithm*. Addison-Wesley, 1996.
- W. Szpankowski. *Average Case Analysis of Algorithms on Sequences*. Wiley, New York, 2001.
- W. Szpankowski. On the height of digital trees and related problems. *Algorithmica*, 6:256–277, 1991.