

# UNIVERSAL ASYMPTOTICS FOR RANDOM TRIES AND PATRICIA TREES

Luc Devroye  
School of Computer Science  
McGill University  
Montreal, Canada H3A 2K6  
luc@cs.mcgill.ca

March 2, 2004

ABSTRACT. We consider random tries and random PATRICIA trees constructed from  $n$  independent strings of symbols drawn from any distribution on any discrete space. We show that many parameters  $Z_n$  of these random structures are universally stable in the sense that  $Z_n/E\{Z_n\}$  tends to one probability. This occurs, for example, when  $Z_n$  is the height, the size, the depth of the last node added, the number of nodes at a given depth (also called the profile), the search time for a partial match, the stack size, or the number of nodes with  $k$  children. These properties are valid without any conditions on the string distributions.

KEYWORDS AND PHRASES. Trie, PATRICIA tree, probabilistic analysis, law of large numbers, concentration inequality, height of a tree.

CR CATEGORIES: 3.74, 5.25, 5.5.

1991 MATHEMATICS SUBJECT CLASSIFICATIONS: 60D05, 68U05.

---

Author' address: School of Computer Science, McGill University, 3480 University Street, Montreal, Canada H3A 2K6. The author' research was sponsored by NSERC Grant A3456 and FCAR Grant 90-ER-0291. FAX number: 1-514-398-3883.

## Introduction

**Tries** are efficient data structures that were initially developed and analyzed by Fredkin (1960) and Knuth (1973). The tries considered here are constructed from  $n$  independent strings  $X_1, \dots, X_n$ , each drawn from  $\prod_{i=1}^{\infty} \Omega_i$ , where  $\Omega_i$ , the  $i$ -th alphabet, is a countable set. By appropriate mapping, we can and do assume that for all  $i$ ,  $\Omega_i = \mathcal{Z}$ . In practice, the alphabets are often  $\{0, 1\}$ , but that won't even be necessary for the results in this paper. Each string  $X_i = (X_{i1}, X_{i2}, \dots)$  defines an infinite path in a tree: from the root, we take the  $X_{i1}$ -st child, then its  $X_{i2}$ -st child, and so forth. The collection of nodes and edges visited by the union of the  $n$  paths is the infinite trie. If the  $X_i$ 's are different, then each infinite path ends with a suffix path that is traversed by that string only. If this suffix path for  $X_i$  starts at node  $u$ , then we may trim it by cutting away everything below node  $u$ . This node becomes the leaf representing  $X_i$ . If this process is repeated for each  $X_i$ , we obtain a finite tree with  $n$  leaves, called the trie. PATRICIA is a space efficient improvement of the classical trie discovered by Morrison (1968) and first studied by Knuth (1973). It is simply obtained by removing from the trie all internal nodes with one child. While it still has  $n$  leaves, each non-leaf (or internal) node has two or more children.

This article attempts to unify the concentration results for most trie parameters, and is a continuation of an earlier paper of the author (2002), where universal concentration results were obtained for the height, size and profile of random PATRICIA trees. A real-valued trie parameter  $Z_n$  is said to be concentrated if  $Z_n/\mathbf{E}\{Z_n\} \rightarrow 1$  in probability. It is universally concentrated if this holds for all string distributions, as long as the strings are independent. If it holds whenever the strings are independent and identically distributed, then we say that  $Z_n$  is universally concentrated for i.i.d. input. These results are obtained by powerful exponential inequalities developed most recently by Boucheron, Lugosi and Massart (2000, 2002), which are related to but more easily applicable than their ancestors, the inequalities of Talagrand (1988, 1989, 1990, 1991a-b, 1993a-b, 1994, 1995, 1996a-b). Ledoux (1996a-b), Azuma (1967) and McDiarmid (1989, 1998). For the practicing computer scientist, they imply that analyzing  $\mathbf{E}\{Z_n\}$  often suffices, as the random variable  $Z_n$  is highly likely to be close to its mean most of the time. In addition, the inequalities give explicit numerical support for that closeness, not hidden behind  $O(\cdot)$  terms, and not requiring often tedious calculations of variances (denoted here by  $\mathbf{V}\{Z_n\}$ ).

## String models

We will refer to a number of string models in this paper. The oldest model is that of the i.i.d. symbols: each symbol is drawn from a symbol distribution  $(p_0, p_1, p_2, \dots)$  over the nonnegative integers. We will call this the i.i.d. symbol model.

Text files have given rise to the Markovian string model, in which the symbols in a string form a Markov chain on the nonnegative integers. This is completely characterized

by the distribution of the first symbol, and the transition probability matrix over  $\mathcal{Z} \times \mathcal{Z}$  (Régnier (1988), Szpankowski (1988), Jacquet and Szpankowski (1991) and Pittel (1985)).

Flajolet and Vallée (1998) and their co-workers have considered various models for storing real numbers. Their base model has as symbols the coefficients in the continued fraction expansion of a random variable  $X$ , drawn from a density on  $[0, 1]$ . We call this the continued fraction trie and PATRICIA tree. Alternatively, Devroye (1984, 1992) takes the  $k$ -ary expansion of  $X$ , and lets the symbols be the digits in the expansion. This is the so-called density model. Just as with the continued fraction model, the symbols are in general dependent. However, in the density model, they are asymptotically uniformly distributed on  $Z_k$ . In the continued fraction model, the distribution on the positive integers is asymptotically fixed as well. Thus, both models should in many cases yield results that are valid for most or all densities of  $X$ .

Clément, Flajolet and Vallée (1998, 2001) and Flajolet and Vallée (1998) consider dynamic sources. Here, the real number  $X$  starts off functional iterations, leading to  $X, f(X), f(f(X)), \dots$ , where  $f : [0, 1] \rightarrow [0, 1]$  is a mapping. There is a second mapping  $s : [0, 1] \rightarrow Z_+$ , the nonnegative integers, and the string symbols are

$$s(X), s(f(X)), s(f(f(X))), \dots$$

For example, by taking  $f(x) = kx \bmod 1$ , and defining  $s(\cdot)$  by partitioning the unit interval equally among  $0, 1, \dots, k - 1$ , we obtain the density model. The continued fraction and Markovian models can be obtained as special cases as well. The first-order parameters of the resulting tries were studied by Flajolet and Vallée (1998) and Clément, Flajolet and Vallée (2001). For random PATRICIA trees, they can be found in the thesis of Bourdon (2002).

Tries are also used as multidimensional data structures. They were first introduced by Orenstein (1982) for database applications. Related ideas had earlier been proposed by Bentley and Burkhard (1976). If  $X$  is a random vector of  $[0, 1]^d$ , then any of the ways of transforming a real number into a string of integers described above can be used. The simplest model uses the  $k$ -ary expansions of each component of  $X$  to generate a string of symbols, each taking values in  $Z_k^d$ . The resulting tries have a possible fanout of  $k^d$ . With  $k = 2$ , they are called quadtries, useful data structures for the compaction of multidimensional (geometric, video) information. In the present paper, parameters such as height, depth, profile and size are dealt with uniformly for all models of tries. The multidimensional tries become interesting in their own right when one considers multivariate operations such as partial match. Puech and Yahia (1985) provide the first analytical study.

## Boucheron-Lugosi-Massart inequality

The following inequalities will be fundamental for the remainder of the paper. Lemma 1 is an almost trivial extension of a similar inequality due to Boucheron, Lugosi and Massart (2000). Its proof is based on logarithmic Sobolev inequalities developed in part by Ledoux (1996a).

LEMMA 1 (BOUCHERON, LUGOSI AND MASSART, 2000). *Let  $\Omega = \mathcal{Z}^n$ . Let  $f \geq 0$  be a function on  $\Omega$ , let  $c \geq 0$  be a constant, and let  $g$  be a real-valued function on  $\mathcal{Z}^{n-1}$  satisfying the following properties for every  $x = (x_1, \dots, x_n) \in \Omega$ :*

$$0 \leq f(x) - g(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \leq 1, \quad 1 \leq i \leq n;$$

$$\sum_{i=1}^n (f(x) - g(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)) \leq f(x) + c.$$

Then for any  $X = (X_1, \dots, X_n)$  with independent components  $X_i \in \mathcal{Z}$ , and all  $t \geq 0$ ,

$$\mathbf{P}\{f(X) \geq \mathbf{E}\{f(X)\} + t\} \leq \exp\left(-\frac{t^2}{2\mathbf{E}\{f(X) + c\} + 2t/3}\right)$$

and

$$\mathbf{P}\{f(X) \leq \mathbf{E}\{f(X)\} - t\} \leq \exp\left(-\frac{t^2}{2\mathbf{E}\{f(X) + c\}}\right).$$

The most outstanding application area for these inequalities are Talagrand's configuration functions. However, as we need to define  $g$  on a space of dimension one less than  $n$ , it is best to reformulate things in terms of "properties". Assume that we have a property  $P$  defined over the union of all finite products  $\mathcal{Z}^k$ . Thus, if  $i_1 < \dots < i_k$ , we have an indicator function that decides whether  $(x_{i_1}, \dots, x_{i_k}) \in \mathcal{Z}^k$  satisfies property  $P$ . We assume that  $P$  is hereditary in the sense that if  $(x_{i_1}, \dots, x_{i_k})$  satisfies  $P$ , then so does any subsequence  $(x_{j_1}, \dots, x_{j_\ell})$  where  $\{j_1, \dots, j_\ell\} \subseteq \{i_1, \dots, i_k\}$ , with the  $j_m$ 's increasing. The configuration function  $f_n(x_1, \dots, x_n)$  gives the size of the largest subsequence of  $x_{i_1}, \dots, x_{i_n}$  satisfying  $P$ . Any subsequence of maximal length satisfying property  $P$  is called a witness. In Lemma 1, we can set  $f(x_1, \dots, x_n) = f_n(x_1, \dots, x_n)$  and  $g(x_1, \dots, x_{n-1}) = f_{n-1}(x_1, \dots, x_{n-1})$ . Clearly, the first condition of Lemma 1 is satisfied, as adding a point to a sequence can only increase the value of the configuration function (so,  $f \geq g$ ), but by not more than one. To verify the second condition, let  $\{x_{i_1}, \dots, x_{i_k}\} \subseteq \{x_1, \dots, x_n\}$  be a witness of the fact that  $f(x_1, \dots, x_n) = k$ . For  $i \leq n$  and  $x_i \notin \{x_{i_1}, \dots, x_{i_k}\}$ , we have  $f(x_1, \dots, x_n) = g(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ , and thus, the difference between  $f$  and  $g$  in the second condition can only be one if  $x_i \in \{x_{i_1}, \dots, x_{i_k}\}$ . Therefore, the sum in that condition is at most  $k = f(x_1, \dots, x_n)$ . Properties  $P$  include being monotonically increasing, being in convex position, and belonging to a given set  $S$ .

We also need some exponential versions of the Efron-Stein inequality. To this end, let  $X_1, \dots, X_n$  be independent random variables in a measurable space, and let  $X'_1, \dots, X'_n$  be an independent copy. Let  $f$  be a measurable mapping, and set

$$Z = f(X_1, \dots, X_n)$$

and

$$Z_i = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n) .$$

We write  $X = (X_1, \dots, X_n)$ . Conditional on  $X$ ,  $Z_i$  is thus a function of  $X'_i$  only. Define

$$V_+ = \mathbf{E} \left\{ \sum_{i=1}^n (Z - Z_i)^2 \mathbf{1}_{[Z > Z_i]} \middle| X \right\}$$

and

$$V_- = \mathbf{E} \left\{ \sum_{i=1}^n (Z - Z_i)^2 \mathbf{1}_{[Z < Z_i]} \middle| X \right\} .$$

The Efron-Stein inequality (Efron and Stein, 1981) states that

$$\mathbf{V}\{Z\} \leq \mathbf{E}\{V_+\} = \mathbf{E}\{V_-\} .$$

By Chebyshev's inequality, we see that if  $Z \geq 0$  is a positive random variable, and  $\epsilon > 0$ ,

$$\mathbf{P}\{|Z - \mathbf{E}\{Z\}| \geq \epsilon \mathbf{E}\{Z\}\} \leq \frac{\mathbf{V}\{Z\}}{\epsilon^2 \mathbf{E}^2\{Z\}} \leq \frac{\mathbf{E}\{V_+\}}{\epsilon^2 \mathbf{E}^2\{Z\}} .$$

In many cases (examples will follow), this ratio tends to zero. However, it does so rather slowly, and the resulting bounds are sometimes unsatisfactory. In this sense, the following inequality is very useful:

LEMMA 2 (BOUCHERON, LUGOSI AND MASSART, 2002). *Assume  $Z \geq 0$ . If  $V_+ \leq aZ + b$  for some nonnegative constants  $a, b$ , then for all  $t > 0$ ,*

$$\mathbf{P}\{Z \geq \mathbf{E}\{Z\} + t\} \leq \exp \left( - \frac{t^2}{4a\mathbf{E}\{Z\} + 4b + 2at} \right) .$$

*If there exists a nondecreasing function  $h$  such that*

$$V_- \leq h(Z) ,$$

*then*

$$\mathbf{P}\{Z \leq \mathbf{E}\{Z\} - t\} \leq \exp \left( - \frac{t^2}{4\mathbf{E}\{h(Z)\}} \right) .$$

Finally, the following inequality is useful whenever only  $V_+$  is easy to bound.

LEMMA 3 (BOUCHERON, LUGOSI AND MASSART, 2002). Assume  $Z \geq 0$  and assume that for all  $i$ ,  $|Z - Z_i| \leq \alpha$ . If  $V_+ \leq \alpha^2 h(Z/\alpha)$  for some nondecreasing function  $h$ , then for all  $0 < t < (e - 1)\alpha \mathbf{E}\{h(Z/\alpha)\}$ ,

$$\mathbf{P}\{Z \leq \mathbf{E}\{Z\} - t\} \leq \exp\left(-\frac{t^2}{4\alpha^2(e-1)\mathbf{E}\{h(Z/\alpha)\}}\right).$$

### Occupancy problems

Consider a very general bin model in which we have  $n$  balls thrown independently into a countable number of bins, where each ball has its own distribution over the bins. Let  $N_1, N_2, \dots$  be the numbers of balls in the bins. Quantities of interest in certain applications include  $M_n = \max_i N_i$ , the maximum number of balls in a single bin, and  $O_n = \sum_i 1_{N_i > 0}$ , the number of occupied bins. If we throw one less ball, then  $M_n$  and  $O_n$  both decrease by at most one. Thus, uniformly over distributions, by the bounded difference inequality (Azuma, 1967; McDiarmid, 1989), we have

$$\mathbf{P}\{|O_n - \mathbf{E}\{O_n\}| \geq t\} \leq 2e^{-t^2/2n}.$$

Also,

$$\mathbf{P}\{|M_n - \mathbf{E}\{M_n\}| \geq t\} \leq 2e^{-t^2/2n}.$$

These results are sometimes unsatisfactory, as  $t$  needs to be at least  $\Omega(\sqrt{n})$  for the inequalities to kick in. Note however that both  $O_n$  and  $M_n$  may be cast in the format of Lemma 1, with  $M_n$  being the configuration function for the hereditary property “belonging to the same bin”, and  $O_n$  being the configuration function for the hereditary property “belonging to different bins”. Thus, by Lemma 1,

$$\mathbf{P}\{O_n \geq \mathbf{E}\{O_n\} + t\} \leq \exp\left(-\frac{t^2}{2\mathbf{E}\{O_n\} + 2t/3}\right), \quad t \geq 0,$$

and

$$\mathbf{P}\{O_n \leq \mathbf{E}\{O_n\} - t\} \leq \exp\left(-\frac{t^2}{2\mathbf{E}\{O_n\}}\right), \quad t \geq 0.$$

Also, for fixed  $t > 0$ , if  $\mathbf{E}\{O_n\} \rightarrow \infty$ ,

$$\mathbf{P}\left\{\left|\frac{O_n - \mathbf{E}\{O_n\}}{\sqrt{\mathbf{E}\{O_n\}}}\right| \geq t\right\} \leq 2 \exp\left(-\frac{t^2}{2 + o(1)}\right).$$

And precisely the same inequalities hold when  $O_n$  is replaced by  $M_n$  throughout. Note that these inequalities are strong enough to imply the following:

$$\frac{O_n}{\mathbf{E}\{O_n\}} \rightarrow 1$$

in probability whenever  $\mathbf{E}\{O_n\} \rightarrow \infty$ , and the result is true over a triangular array of bin distributions. Also, we have

$$\frac{M_n}{\mathbf{E}\{M_n\}} \rightarrow 1$$

in probability whenever  $\mathbf{E}\{M_n\} \rightarrow \infty$ .

In data structures, these results are relevant for hashing with chaining with equal or unequal probabilities. The maximal chain length satisfies the law of large numbers regardless of how the table size changes with  $n$ . For  $M_n$ , if the number of bins equals the number of balls, then  $M_n \sim \log n / \log \log n$  if each bin has equal probability of receiving a ball. The inequalities at the top of the section would not allow one to obtain a law of large numbers. However, Lemma 1, as shown above, suffices to obtain it. See Gonnet (1981), Devroye (1985), or Knuth (1973) for more on the maximum chain length.

It is interesting to note that problems on tries can often be regarded as problems on bins. In fact, we consider levels of nested bins, with the 0-th level consisting of one bin, corresponding to the root, the first level having one bin for each child of the root, and so forth. A string  $X_i = (X_{i1}, X_{i2}, \dots)$  thus drops a ball in the bin characterized by  $X_{i1}$  of the level 1 collection, a ball in the bin characterized by  $(X_{i1}, X_{i2})$  of the level 2 collection, and so forth. We call bin  $(X_{i1}, \dots, X_{i,j+1})$  a child bin of  $(X_{i1}, \dots, X_{i,j})$ . In a trie, a bin that receives at least two balls corresponds to an internal node. The number of nodes is  $n$  (the number of leaves) plus the number of internal nodes, which is a random variable. In a PATRICIA tree, a bin with at least two balls, whose parent bin has more balls, corresponds to an internal node. These observations permit us to make observations about the size and the profile of random tries and PATRICIA trees.

### Size of a PATRICIA tree

Let  $S_n$  be the number of internal nodes, and let  $T_n = S_n + n$  be the total number of nodes in a PATRICIA tree for  $n$  strings. Note that for binary PATRICIA trees,  $S_n = n - 1$ , so only non-binary trees have random sizes. Adding a string increases  $T_n$  by one and  $S_n$  by one or zero. Thus, if the strings are independent (but not necessarily identically distributed), by the bounded difference inequality (McDiarmid, 1989),

$$\mathbf{P}\{|S_n - \mathbf{E}\{S_n\}| \geq t\} = \mathbf{P}\{|T_n - \mathbf{E}\{T_n\}| \geq t\} \leq 2 \exp\left(-\frac{t^2}{2n}\right).$$

The fanout and string distributions do not figure in the bound. We immediately have

$$\frac{T_n}{\mathbf{E}\{T_n\}} \rightarrow 1$$

almost surely (as  $T_n \geq n$ ), and

$$\frac{S_n}{\mathbf{E}\{S_n\}} \rightarrow 1$$

in probability whenever  $\mathbf{E}\{S_n\}/\sqrt{n} \rightarrow \infty$  (which is satisfied, for example, if the strings consist of independent identically distributed symbols, or when the tree is of bounded fanout). Even though these results do not require Lemma 1, they appear to be new. Still, we would like to offer a stronger result.

Let  $Z = S_n$  be the parameter of interest. Observe that removing a string  $X_i$  causes

$Z$  to shrink by at most one node, so that  $|Z - Z_i| \leq 1$ . Furthermore,

$$\sum_i (Z - Z_i)^2 \mathbf{1}_{[Z > Z_i]} \leq \sum_i \mathbf{1}_{[Z > Z_i]} \leq 2Z$$

as each internal node can be removed only when it has two balls and one of its two balls are removed. It cannot be removed when it has more than two balls. Thus,  $V_+ \leq 2Z$ . We are in a position to apply Lemma 2 with  $(a, b) = (2, 0)$ , and obtain

$$\mathbf{P}\{S_n \geq \mathbf{E}\{S_n\} + t\} \leq \exp\left(-\frac{t^2}{8\mathbf{E}\{S_n\} + 4t}\right).$$

Furthermore, we can apply Lemma 3 with  $h(u) = 2u$ ,  $\alpha = 1$  and obtain

$$\mathbf{P}\{S_n \leq \mathbf{E}\{S_n\} - t\} \leq \exp\left(-\frac{t^2}{8(e-1)\mathbf{E}\{S_n\}}\right), \quad 0 < t < 2(e-1)\mathbf{E}\{S_n\}.$$

By setting  $t = \epsilon \mathbf{E}\{S_n\}$ , we conclude immediately that

$$\frac{S_n}{\mathbf{E}\{S_n\}} \rightarrow 1 \text{ in probability}$$

whenever  $\mathbf{E}\{S_n\} \rightarrow \infty$ . Only pathological examples do not satisfy this condition. For example, this would occur if  $X_i = (i, i, i, \dots)$ , so that the PATRICIA tree consists of one root and  $n$  child leaves, and thus,  $S_n = 1$ .

When the strings are independent and identically distributed, and the symbols in the strings are independent and identically distributed as well, the size  $S_n$  of a random trie has been studied by Régnier and Jacquet (1989) and Jacquet and Régnier (1989). The expected value  $\mathbf{E}\{S_n\}$  is close to  $n/\mathcal{H}$ , where  $\mathcal{H}$  is the entropy of a random symbol. (If a symbol has distribution  $(p_j)_{j \geq 0}$ , then the entropy  $\mathcal{H}_k$  of order  $k$  is  $\sum_j p_j \log^k(1/p_j)$ . The entropy is the entropy of order one.) There are string distributions that cause  $\mathbf{E}\{S_2\} = \infty$ : just let each string be all zeroes, with ones added in at places  $W_1, W_1 + W_2, \dots$ , and the  $W_i$ 's are i.i.d. with the same distribution on the positive integers. Then  $S_2 \geq \min(W_1, W_1')$ , the minimum of two independent copies of  $W_1$ . Thus,  $\mathbf{P}\{\min(W_1, W_1') \geq t\} = (\mathbf{P}\{W_1 \geq t\})^2$  and so, whenever  $\mathbf{P}\{W_1 \geq t\} \geq C/\sqrt{t}$ , for some positive  $C$ , we have  $\mathbf{E}\{S_2\} = \infty$ .

Rais, Jacquet and Szpankowski (1993) studied  $S_n$  for PATRICIA trees when the symbols are independent. Bourdon (2001, 2002) studied  $S_n$  for PATRICIA trees under a general string distribution model. It is noted that in the former case,  $\mathbf{E}\{S_n\}$  is about  $(1 - \mathcal{H}_2/(2\mathcal{H}^2))n/\mathcal{H}$ . With more general definitions of entropy for dependent sequences of symbols (as in Markovian sources), the same result was obtained by Bourdon. For the continued fraction trie, Bourdon showed that  $\mathbf{E}\{S_n\}/n$  is about

$$\frac{6 \log 2}{\pi^2}$$

(this is the inverse of Lévy's constant), and for the random PATRICIA tree, he showed that it is about equal to the same constant multiplied by

$$\left(1 - \frac{12\gamma \log 2}{\pi^2} - \frac{9 \log^2 2}{\pi^2} + \frac{72\zeta'(2) \log 2}{\pi^4} + 1/2\right) \approx 0.87.$$



Here  $\gamma$  is Euler’s constant, and  $\zeta$  is Riemann’s zeta function. The study of variances and second order properties of  $S_n$  for PATRICIA trees was left as a major open problem in the thesis of Bourdon (2002). For symbol distributions that are “periodic”, the expected value of  $S_n/n$  oscillates forever for random tries and PATRICIA trees, so the limiting constants mentioned above are not true “limits”. However, as we showed in this paper, even in those cases,  $S_n/\mathbb{E}\{S_n\} \rightarrow 1$  in probability. We did not have to worry about the oscillation phenomenon. We observe in fact, that for all models, without exception,

$$\mathbb{V}\{S_n\} \leq \mathbb{E}\{V_+\} \leq 2\mathbb{E}\{S_n\} .$$

### Multidimensional tries and the partial match operation

We can define a general query in a trie or PATRICIA tree as follows: it is defined by a sequence of sets  $S_1, S_2, S_3, \dots$  of symbols, and asks for all the leaves (if any) of the trie or PATRICIA tree that are in the infinite query tree as well, where the query tree consists of a root, all nodes at level one indexed by  $s_1 \in S_1$ , all nodes at level 2 indexed by  $(s_1, s_2) \in S_1 \times S_2$ , and so forth. Thus, the query tree can be viewed as an infinite trie. The time needed to collect all the leaves that need reporting is proportional to the number  $N$  of nodes in the intersection of the query tree and the original trie or PATRICIA tree. In the case of a PATRICIA, care must be taken to identify nodes in the tree not by their position but by their index. Some examples follow that show that most parameters in a random trie can be viewed as special cases of  $N$ .

EXAMPLE 1. If  $S_1 = S_2 = \dots$  is the full symbol set, then we must report all leaves, and the time  $N$  in that case is the size of the random trie.

EXAMPLE 2. If  $S_i = \{y_i\}$  for all  $i$ , then the query tree is the infinite path  $(y_1, y_2, y_3, \dots)$ , and thus,  $N$  is the size of the path followed when we search for the string  $(y_1, y_2, y_3, \dots)$  in the trie or PATRICIA tree. This is equal to the depth of the unique leaf on that path plus one. Thus, search times for individual strings are obtained as special cases.

EXAMPLE 3. In the multivariate random trie and random PATRICIA tree, with a fanout of  $k^d$  (see the introduction for the notation), we may describe each symbol as an element of  $Z_k^d$ . A partial match query can be regarded as a general query in which  $S_i = S_{i1} \times \dots \times S_{id}$ ,  $S_{ij} = Z_k$  for all  $j \in J \subseteq \{1, 2, \dots, d\}$ , the so-called set of wild cards. For  $j \notin J$ , we have  $S_{ij} = \{y_{ij}\}$ , a singleton set. The partial match is thus determined by all the values  $y_{ij}, i \geq 1, j \notin J$ . Usually, the interpretation is easy in term of a vector  $y \in [0, 1]^d$ , whose

$j$ -th component has a  $k$ -ary expansion  $(y_{1j}, y_{2j}, \dots)$ . Thus, the pair  $(y, J)$  uniquely defines a partial match query. We may thus write  $N(y, J)$  to denote the size of the intersection of our two trees for the query given by  $(y, J)$ . If  $J$  is the full set, then the query tree is the full tree, and we are reduced to example 1. If  $J$  is empty, then the query tree reduces to a path, as in example 2. The only interesting cases appear when  $0 < |J| < d$ , and in that case, we speak of a proper partial match. Few results are known, and they all relate to random tries with  $k = 2$ , when each of the strings in the random trie is generated by an independently drawn random vector uniformly distributed on  $[0, 1]^d$ . In that case, all symbols in all strings are independent and uniformly distributed on  $Z_2^d$ . Consider a proper partial match. Flajolet and Puech (1986) showed that

$$\mathbf{E}\{N(y, J)\} = \tau(\log_2 n) n^{|J|/d} + o\left(n^{|J|/d}\right)$$

where  $\tau$  is a continuous positive periodic function. Kirschenhofer, Prodinger and Szpankowski (1993) showed that for  $k = 2$ ,  $d = 2$ ,  $|J| = 1$ ,

$$\mathbf{V}\{N(y, J)\} \sim \tau(\log_2 n) \sqrt{n}$$

where  $\tau$  is again a continuous positive periodic function. This was generalized to  $d > 2$  by Schachinger (1995). In 2000, Schachinger proved that

$$\frac{N(y, J)}{\mathbf{E}\{N(y, J)\}} \rightarrow 1$$

in probability as  $n \rightarrow \infty$ , and

$$\frac{N(y, J) - \mathbf{E}\{N(y, J)\}}{\sqrt{\mathbf{V}\{N(y, J)\}}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1) ,$$

where  $\xrightarrow{\mathcal{L}}$  denotes convergence in distribution. Using concentration inequalities, Devroye and Zamora-Cura (2002) were able to show that

$$\frac{\sup_y N(y, J)}{\inf_y N(y, J)} \rightarrow 1$$

in probability, which shows a remarkable stability of all partial match query times: every partial match query time is with high probability close to  $\mathbf{E}\{N(y_0, J)\}$ , where  $y_0$  is the vector of all zeroes. The asymmetric Bernoulli model refers to strings with independent identically distributed symbols, and each symbol consists of  $d$  independent components, with each component drawn from some non-uniform distribution of  $Z_2$ . For this model, some asymptotics for  $\mathbf{E}\{N(y, J)\}$  are obtained by Kirschenhofer, Prodinger and Szpankowski (1993), and, with  $y$  replaced by a random  $Y$ , by Schachinger (2000). In the latter paper, it is shown that  $N(y, J)/\mathbf{E}\{N(Y, J)\} \xrightarrow{\mathcal{L}} Z_p$  under an idealized partial match model that does not correspond to our definition, where the distribution of  $Z_p$  depends upon the probability  $p$  only. More recently, Schachinger has studied the asymptotic behavior of the ratio  $\log N(y, J)/\mathbf{E}\{\log N(y, J)\}$ .

EXAMPLE 4. If we take all  $S_i, 1 \leq i \leq k$  equal to the full set of symbols, and all  $S_i, i > k$  empty, then we obtain the cumulative size of the tree up to level  $k$ . It is also called the

cumulative profile. In fact, the profile would be obtained by considering all the nodes at one level. Strictly speaking, the profile does not follow the subtree analogy set up in this section, so we will deal with it below.

So, we consider  $N$  for a general query, freed from the multidimensional confines. This  $N$  is our random variable to which we want to apply some concentration inequalities. When a string is removed, in the PATRICIA version,  $N$  decreases by at most two, as one leaf disappears and possibly one internal node. Thus, denoting by  $N_i$  the size of the intersection tree when string  $X_i$  is removed but an independent copy  $X'_i$  is added, we see that  $|N - N_i| \leq 2$ . Furthermore,

$$\sum_i (N - N_i)^2 \mathbf{1}_{[N > N_i]} \leq 4 \sum_i \mathbf{1}_{[N > N_i]} \leq 8N$$

as each internal node can be removed only when it has two balls and one of its two balls is removed. It cannot be removed when it has more than two balls. Thus,  $V_+ \leq 8N$ . We are in a position to apply Lemma 2 with  $(a, b) = (8, 0)$ , and obtain

$$\mathbf{P}\{N \geq \mathbf{E}\{N\} + t\} \leq \exp\left(-\frac{t^2}{32\mathbf{E}\{N\} + 16t}\right).$$

Furthermore, we can apply Lemma 3 with  $h(u) = 4u$ ,  $\alpha = 2$  and obtain

$$\mathbf{P}\{N \leq \mathbf{E}\{N\} - t\} \leq \exp\left(-\frac{t^2}{32(e-1)\mathbf{E}\{N\}}\right), \quad 0 < t < 8(e-1)\mathbf{E}\{N\}.$$

By setting  $t = \epsilon\mathbf{E}\{N\}$ , we conclude immediately that

$$\frac{N}{\mathbf{E}\{N\}} \rightarrow 1 \text{ in probability}$$

whenever  $\mathbf{E}\{N\} \rightarrow \infty$ .

Finally, the Efron-Stein inequality implies that for all models, without exception, and for any kind of general query,

$$\mathbf{V}\{N\} \leq \mathbf{E}\{V_+\} \leq 8\mathbf{E}\{N\}.$$

It is interesting to note that the addition of one string to a trie can cause a trie to grow by many nodes, and that this growth cannot be a priori bounded. The concentration inequalities given here can thus not be used directly. However, if  $X_i$  is removed and replaced by string  $X'_i$ , we have

$$\sum_i (N - N_i)^2 \mathbf{1}_{[N > N_i]} \leq (1 + H_n)^2 \sum_i \mathbf{1}_{[N > N_i]} \leq 2(1 + H_n)^2 N$$

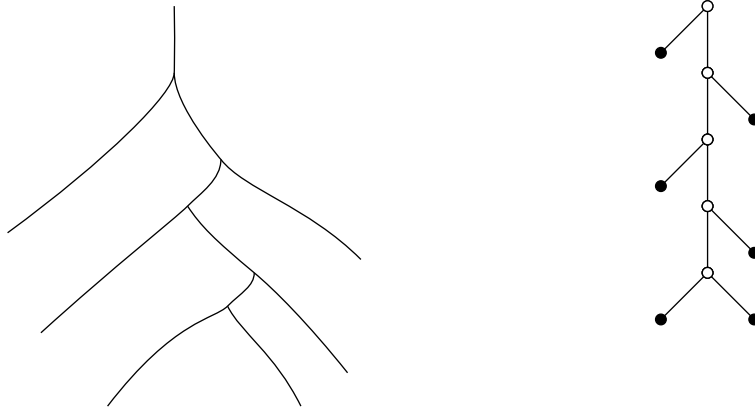
where  $H_n$  is height of the random trie, as  $|N - N_i| \leq 1 + H_n$ . Thus,  $V_+ \leq 2(1 + H_n)^2 N$ . In most, but not all, models,  $H_n$  is of the order of  $\log n$ , so that one can expect to make good use of this if one has further information on the height  $H_n$ , and this requires, unfortunately, some information on the string distributions. We would like to draw the attention though to an inequality of Boucheron, Lugosi and Massart (2002) that deals precisely with the case that  $V_+ \leq WZ$ , where  $Z$  is the random variable of interest, and  $W$  is another random

variable, whose tail can be bounded sufficiently tightly. This way, one can obtain good concentration inequalities, although not universal ones, for  $N$  in random tries.

The second way in which  $N$  can be dealt with in random tries is by decomposing  $N = \sum_{\ell} N^{(\ell)}$ , where  $N^{(\ell)}$  is the number of nodes in the intersection with the query tree that are at level  $\ell$ . As  $N^{(\ell)}$  is easy to deal with, one can obtain a concentration inequality for it (see the section on general profiles). It is also true that in many models,  $N$  is nearly equal to the sum of  $N^{(\ell)}$  with  $\ell$  ranging over just a few levels. In those cases, concentration for  $N$  follows. An example of this approach is worked out by Devroye and Zamora-Cura (2002) for partial match queries in multidimensional tries in which all symbols are i.i.d. and uniformly distributed on  $Z_2^d$ .

### Height of a PATRICIA tree

Given are  $n$  independent infinite strings  $X_1, \dots, X_n$  (if they are not infinite, pad them by some designated character, repeated infinitely often), each drawn from a distribution on  $\mathcal{Z}$ . The height of the PATRICIA tree is denoted by  $H_n$ . If (deterministic) strings  $x_1, \dots, x_k$  induce a PATRICIA tree of height  $k - 1$ , then the PATRICIA tree can have only one configuration, namely, it consists of a chain of length  $k - 1$  from the root on down, with every node of this chain receiving one leaf, except the furthest node, which receives two leaves. We say that such a collection of strings has the PATRICIA property. This property is clearly hereditary, and  $H_n + 1$  is thus a configuration function.



Six strings with the PATRICIA property. Each (black) leaf represents a contracted infinite string. The height is five.

We have

$$\mathbf{P}\{H_n \geq \mathbf{E}\{H_n\} + t\} \leq \exp\left(-\frac{t^2}{2\mathbf{E}\{(H_n + 1)\} + 2t/3}\right), \quad t \geq 0,$$

and

$$\mathbb{P}\{H_n \leq \mathbb{E}\{H_n\} - t\} \leq \exp\left(-\frac{t^2}{2\mathbb{E}\{(H_n + 1)\}}\right), \quad t \geq 0.$$

We stress that the individual strings may have any distribution. The symbols themselves need not be independent or identically distributed. And the strings need not be identically distributed. All PATRICIA trees, without exception, are thus stable and well-behaved. For any PATRICIA tree constructed by using  $n$  independent strings, if  $\lim_{n \rightarrow \infty} \mathbb{E}\{H_n\} = \infty$ , then

$$\frac{H_n}{\mathbb{E}\{H_n\}} \rightarrow 1$$

in probability as  $n \rightarrow \infty$ , and

$$\frac{H_n - \mathbb{E}\{H_n\}}{\sqrt{\mathbb{E}\{H_n\}}} = O(1)$$

in probability in this sense: for fixed  $t > 0$ ,

$$\mathbb{P}\left\{\left|\frac{H_n - \mathbb{E}\{H_n\}}{\sqrt{\mathbb{E}\{H_n\}}}\right| \geq t\right\} \leq 2 \exp\left(-\frac{t^2}{2 + o(1)}\right).$$

The last inequality remains valid whenever  $0 < t = o(\mathbb{E}\{H_n\})$ .

THE CONDITION ON  $\mathbb{E}\{H_n\}$ . In PATRICIA trees of bounded degree, it is clear that  $\mathbb{E}\{H_n\} \rightarrow \infty$ . In unbounded degree trees, this is also true provided that the strings are identically distributed and the probability of two identical strings is zero. However, without the identical distribution constraint, PATRICIA trees may have  $H_n = 1$  for all  $n$ : just let the  $i$ -th string be  $(i, 0, 0, 0, \dots)$ .

BIBLIOGRAPHIC REMARK: HEIGHT OF PATRICIA TREES. All parameters of a PATRICIA tree such as  $H_n$  improve over those of the associated trie: for the uniform trie model, Pittel (1985) has shown that  $H_n / \log_2 n \rightarrow 1$  almost surely, which constitutes a 50% improvement over the trie. For other properties, see Knuth (1973), Flajolet and Sedgewick (1986), Kirschenhofer and Prodinger (1986) and Szpankowski (1990, 1991). Pittel and Rubin (1990), Pittel (1991) and Devroye (1992) showed that

$$\frac{H_n - \log_2 n}{\sqrt{2 \log_2 n}} \rightarrow 1 \quad \text{almost surely.}$$

More refined results for general multi-branching PATRICIA trees and tries are given by Szpankowski and Knessl (2000). For the independent symbol trie model with symbol probabilities  $p_j$ , we have  $\mathbb{E}\{H_n\} \sim c \log n$ , where  $c = 2 / \log_2(1 / \sum_j p_j^2)$ .

## Depth along a given path in a PATRICIA tree

Consider a string  $x$  that defines an infinite path in a trie. We define the depth of the path  $x$ , denoted by  $D_n(x)$  in the PATRICIA tree as the depth (distance to the root) of the leaf that corresponds to  $x$  in the PATRICIA tree for  $X_1, \dots, X_n, x$ . We say that strings  $x_1, \dots, x_k$  have the  $x$ -property if the prefixes  $x \cap x_1, \dots, x \cap x_k$  are strictly nested. That is, there is a reordering  $x'_1, \dots, x'_k$  of the strings such that the common prefix of  $x'_1$  and  $x$  is strictly contained in that of  $x'_2$  and  $x$ , and so forth. In that case, the distance of the leaf of  $x$  from the root of the PATRICIA tree for  $x_1, \dots, x_k, x$  is precisely  $k$ . The function  $D_n(x) = f(x_1, \dots, x_n)$  that describes the length of the longest subset of  $x_1, \dots, x_n$  with the  $x$ -property is clearly a configuration function, to which Lemma 1 may be applied. Thus, we conclude as in the previous section:

For any PATRICIA tree constructed by using  $n$  independent strings, if  $x$  is a string such that  $\lim_{n \rightarrow \infty} \mathbf{E}\{D_n(x)\} = \infty$ , then

$$\frac{D_n(x)}{\mathbf{E}\{D_n(x)\}} \rightarrow 1$$

in probability as  $n \rightarrow \infty$ , and

$$\frac{D_n(x) - \mathbf{E}\{D_n(x)\}}{\sqrt{\mathbf{E}\{D_n(x)\}}} = O(1)$$

in probability in this sense: for fixed  $t > 0$ ,

$$\mathbf{P} \left\{ \left| \frac{D_n(x) - \mathbf{E}\{D_n(x)\}}{\sqrt{\mathbf{E}\{D_n(x)\}}} \right| \geq t \right\} \leq 2 \exp \left( -\frac{t^2}{2 + o(1)} \right).$$

## General profiles

This section is about the profile in random tries. The profile of a trie is the sequence  $(P_1, P_2, P_3, \dots)$ , where  $P_\ell$  is the number of nodes at level  $\ell$  (at distance  $\ell$  from the root). This is in fact the number of nodes that have either at least leaves in their subtrees or correspond to leaves. Let  $P_{\ell i}$  denote the value when string  $X_i$  is deleted and string  $X'_i$  is added. It is easy to see that  $|P_\ell - P_{\ell i}| \leq 1$ . Furthermore,

$$\sum_i (P_\ell - P_{\ell i})^2 \mathbf{1}_{[P_\ell > P_{\ell i}]} \leq \sum_i \mathbf{1}_{[P_\ell > P_{\ell i}]} \leq 2P_\ell$$

because only internal nodes with two balls can be deleted when  $X_i$  is deleted from the list of strings. Summed over all  $i$ , this can happen just twice. So, we have a situation exactly like that for the size  $S_n$  of a PATRICIA tree. In particular,  $V_+ \leq 2P_\ell$ . We apply Lemma 2 with  $(a, b) = (2, 0)$ :

$$\mathbf{P}\{P_\ell \geq \mathbf{E}\{P_\ell\} + t\} \leq \exp \left( -\frac{t^2}{8\mathbf{E}\{P_\ell\} + 4t} \right).$$

Furthermore, we can apply Lemma 3 with  $h(u) = 2u$ ,  $\alpha = 1$  and obtain

$$\mathbf{P}\{P_\ell \leq \mathbf{E}\{P_\ell\} - t\} \leq \exp\left(-\frac{t^2}{8(e-1)\mathbf{E}\{P_\ell\}}\right), \quad 0 < t < 2(e-1)\mathbf{E}\{P_\ell\}.$$

By setting  $t = \epsilon\mathbf{E}\{P_\ell\}$ , we conclude immediately that

$$\frac{P_\ell}{\mathbf{E}\{P_\ell\}} \rightarrow 1 \text{ in probability}$$

whenever  $\mathbf{E}\{P_\ell\} \rightarrow \infty$ , regardless of whether  $\ell$  is fixed or varies with  $n$ . Finally,

$$\mathbf{V}\{P_\ell\} \leq 2\mathbf{E}\{P_\ell\}.$$

All of the above remains valid for the general profile, which for each level counts the number of trie nodes (internal or leaves) that belong to certain set  $S$  on the space of all indices.

### Height of a random trie: lower bound

Note the duality

$$[H_n \geq \ell] = [P_\ell > 0].$$

Thus,

$$\mathbf{P}\{H_n < \ell\} = \mathbf{P}\{P_\ell = 0\} = \mathbf{P}\{P_\ell - \mathbf{E}\{P_\ell\} \leq -\mathbf{E}\{P_\ell\}\} \leq \exp\left(-\frac{\mathbf{E}\{P_\ell\}}{8(e-1)}\right).$$

This inequality is valid without any conditions. For random tries, we have a simple universal lower bound on the tail of the height  $H_n$  in terms of the first moment of  $P_\ell$  only. Since we also have

$$\mathbf{P}\{H_n \geq \ell\} = \mathbf{P}\{P_\ell \geq 1\} \leq \mathbf{E}\{P_\ell\},$$

we can conclude that if  $\ell = \ell(n)$  is a function of  $n$ , then

$$\lim_{n \rightarrow \infty} \mathbf{P}\{H_n \geq \ell(n)\} = \begin{cases} 0 & \text{if } \lim_{n \rightarrow \infty} \mathbf{E}\{P_{\ell(n)}\} = 0, \\ 1 & \text{if } \lim_{n \rightarrow \infty} \mathbf{E}\{P_{\ell(n)}\} = \infty. \end{cases}$$

First-order behavior of the height of all random tries reduces thus simply to the study of the expected profile.

**BIBLIOGRAPHIC REMARK: HEIGHT OF RANDOM TRIES.** The asymptotic behavior of tries under the uniform trie model is well-known. For example, it is known that

$$H_n / \log_2 n \rightarrow 2 \text{ almost surely.}$$

The limit law of  $H_n$  was obtained in Devroye (1984), and laws of the iterated logarithm for the difference  $H_n - 2\log_2 n$  can be found in Devroye (1990). The height for other models was studied by Régnier (1981), Mendelson (1982), Flajolet and Steyaert (1982), Flajolet (1983), Devroye (1984), Pittel (1985, 1986), and Szpankowski (1988, 1989). For the depth of

a node, see e.g., Pittel (1986), Jacquet and Régnier (1986), Flajolet and Sedgewick (1986), Kirschenhofer and Prodinger (1986), and Szpankowski (1988).

### Height of a random trie: upper bound

It is possible to deal with upper bounds for the tail of  $H_n$  in random tries in some way, despite the instability of this parameter for some string distributions. The route suggested here is based on Talagrand's q-points inequality and a subadditivity argument due to van der Vaart and Wellner (1996). Assume that random variables  $X_i \in \Omega$  are given that are independent and identically distributed. Let  $\Omega^n$  be the product space.

LEMMA 4 (VAN DER VAART AND WELLNER, 1996). *Assume that  $\Omega$  is Euclidean space, and that  $f_n : \Omega^n \rightarrow \mathbf{R}$  is a permutation-symmetric function, satisfying the following monotonicity and subadditivity conditions:*

$$\begin{aligned} f_n(x) &\leq f_{n+m}(x, y), x \in \Omega^n, y \in \Omega^m ; \\ f_{n+m}(x, y) &\leq f_n(x) + f_m(y), x \in \Omega^n, y \in \Omega^m . \end{aligned}$$

Assume furthermore that  $0 \leq f_n \leq n$ . Let  $X = (X_1, \dots, X_n)$  have i.i.d. coordinates in  $\Omega^n$ . Then

$$\mathbf{P}\{f_n(X_1, \dots, X_n) \geq t\} \leq \exp\left(-\frac{t}{2} \log\left(\frac{t}{12 \max(1, \mathbf{E}\{f_n(X_1, \dots, X_n)\})}\right)\right)$$

for all  $t > 0$ ,  $n \geq 1$ .

Consider a random trie based upon  $n$  i.i.d. strings  $X_1, \dots, X_n$  with string symbols from an arbitrary alphabet, and let the string distribution be arbitrary. Then the height  $H_n$  satisfies all the required conditions, except possibly the requirement that  $H_n \leq n$ . But  $\min(H_n, n)$  is fully compliant. The Euclidean space requirement is fulfilled as we can study  $\min(H_n, n)$  by just considering string collections truncated to their length  $n$  prefixes. Set  $h_n = \mathbf{E}\{\min(H_n, n)\}$ . Note that for  $n \geq 2$ ,  $h_n \geq 1$ . Then we have for  $n \geq 2$ ,  $0 \leq t \leq n$ ,

$$\mathbf{P}\{H_n \geq t\} = \mathbf{P}\{\min(H_n, n) \geq t\} \leq \exp\left(-\frac{t}{2} \log\left(\frac{t}{12h_n}\right)\right) .$$

Thus, if  $h_n \rightarrow \infty$ , yet  $h_n = o(n)$ , as is usually the case, we see that for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbf{P}\{H_n \geq (12 + \epsilon)h_n\} = 0$$

in all generality. It is difficult to imagine that one can state such universal results without access to the powerful machinery provided by Talagrand's inequalities.



## External path length of a PATRICIA tree

The external path length  $I_n$  of a PATRICIA tree or trie is the sum of the distances from the leaves to the root. In the case of a trie, it is an upper bound for the time needed to construct the trie. Upper bounds for the tail of  $I_n$  are thus useful. When the  $i$ -th string is deleted in a PATRICIA tree, then  $I_n$  decreases by  $D_i$ , the distance from the  $i$ -th leaf to the root, and, in case one internal is removed in the process, by the number of leaves,  $i$  excepted, in the subtree of the node that disappeared (as all those leaves decrease their depth by one). That number of such leaves is denoted by  $N_i$ . Note that  $N_i = 0$  if no leaf changes depth. Now add new independent string  $X_i'$ . Denote by  $I_{ni}$  the new external path length. We thus have

$$\sum_i (I_n - I_{ni})^2 1_{\{I_n > I_{ni}\}} \leq \sum_i (D_i + N_i)^2 \leq n \sum_i (D_i + N_i) \leq 2nI_n ,$$

where we used the fact that  $\sum_i N_i$  is bounded from above by  $I_n$ , as each leaf contributes at most its depth towards the total sum. Thus, we are in a position to apply Lemma 2 with  $(a, b) = (2n, 0)$ . We see that for all  $t > 0$ ,

$$\mathbf{P}\{I_n \geq \mathbf{E}\{I_n\} + t\} \leq \exp\left(-\frac{t^2}{8n\mathbf{E}\{I_n\} + 4nt}\right).$$

In particular, for  $s > 0$ ,

$$\mathbf{P}\{I_n \geq (1 + s)\mathbf{E}\{I_n\}\} \leq \exp\left(-\frac{s^2\mathbf{E}\{I_n\}}{8n + 4ns}\right).$$

This probability tends to zero whenever  $\mathbf{E}\{I_n\}/n \rightarrow \infty$ , a condition that is satisfied for all models with a finite symbol alphabet, and nearly all models with infinite symbol alphabet as well. In fact  $\mathbf{E}\{I_n\} = \sum_i \mathbf{E}\{D_i\} = n\mathbf{E}\{D_1\}$ , so that the condition holds if and only if  $\mathbf{E}\{D_1\} \rightarrow \infty$ .

## Other parameters

We can deal with other parameters of tries and PATRICIA trees with equal ease. Easiest among these is the Horton-Strahler number, which measures the minimum number of registers needed for evaluating an expression when the tree is considered as an expression tree. This quantity is always bounded by  $O(\log_2 n)$  and can change by at most one when any string is deleted or added. It is thus stable in a strong sense, and concentration inequalities are easily derived. In the context of tries, this quantity is however less important, as most models of random tries do not qualify as good models of expression trees.

The stack size of a tree is the maximal stack size needed when traversing a tree in pre-order, while the last subtree (the last trie symbol) is not put on the stack. It is bounded by the height of the tree, and is analyzed for a number of random trie models by Bourdon, Nebel and Vallée (2001). For the PATRICIA tree, this parameter can be dealt with using the configuration function method used for the height. The details are left to the reader.

If  $N$  is the number of nodes in a PATRICIA tree with  $k$  children, then adding a new string can increase  $N$  by at most one. The Azuma-McDiarmid inequality thus immediately implies

$$\mathbb{P}\{|N - \mathbb{E}\{N\}| \geq t\} \leq 2e^{-t^2/2n} .$$

This should suffice for most situations. In this case,  $V_+$  is bounded by the number of nodes with  $k - 1$  children, and thus, a different argument than that provided by Lemmas 2 and 3 is needed, if one wants to improve on the Azuma-McDiarmid inequality.

Our analysis still requires the strings to be independent. It is an interesting challenge to deal with dependent strings, such as those occurring in suffix tries and suffix trees. In those cases, we will invariably have to place conditions on the string symbols themselves, as most concentration inequalities, at their core, require independence. This will be dealt with elsewhere.

## Acknowledgment

The author would like to thank both referees for their help.

## References

- D. Aldous and P. Shields, “A diffusion limit for a class of randomly-growing binary trees,” *Probability Theory and Related Fields*, vol. 79, pp. 509–542, 1988.
- K. Azuma, “Weighted sums of certain dependent random variables,” *Tohoku Mathematical Journal*, vol. 37, pp. 357–367, 1967.
- J. L. Bentley and W. A. Burkhard, “Heuristics for partial-match retrieval in database design,” *Information Processing Letters*, vol. 4(5), pp. 132–135, 1976.
- S. Boucheron, G. Lugosi, and P. Massart, “A sharp concentration inequality with applications in random combinatorics and learning,” *Random Structures and Algorithms*, vol. 16, pp. 277–292, 2000.
- S. Boucheron, G. Lugosi, and P. Massart, “Concentration inequalities using the entropy method,” *The Annals of Probability*, 2002. To appear.
- J. Bourdon, “Size and path length of Patricia tries: dynamical sources context,” *Random Structures and Algorithms*, vol. 19, pp. 289–315, 2001.
- J. Bourdon, M. Nebel, and B. Vallée, “On the stack size of general tries,” *Theoretical Informatics and Applications*, vol. 35, pp. 163–185, 2001.
- J. Bourdon, “Analyse dynamique d’algorithmes: exemples en arithmétique et en théorie de l’information,” Thèse de doctorat, Université de Caen Basse-Normandie, 2002.

- J. Clément, P. Flajolet, and B. Vallée, “The analysis of hybrid trie structures,” in: *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (San Francisco, CA, 1998)*, pp. 531–539, ACM, New York, 1998.
- J. Clément, P. Flajolet, and B. Vallée, “Dynamical sources in information theory: a general analysis of trie structures,” *Algorithmica*, vol. 29, pp. 307–369, 2001.
- L. Devroye, “A probabilistic analysis of the height of tries and of the complexity of triesort,” *Acta Informatica*, vol. 21, pp. 229–237, 1984.
- L. Devroye, “The expected length of the longest probe sequence when the distribution is not uniform,” *Journal of Algorithms*, vol. 6, pp. 1–9, 1985.
- L. Devroye, “A note on the probabilistic analysis of PATRICIA trees,” *Random Structures and Algorithms*, vol. 3, pp. 203–214, 1992.
- L. Devroye, “A study of trie-like structures under the density model,” *Annals of Applied Probability*, vol. 2, pp. 402–434, 1992.
- L. Devroye and C. Zamora-Cura, “Expected worst-case time for partial match in random quadtries,” *Discrete Applied Mathematics*, 2002. To appear.
- L. Devroye, “Laws of large numbers and tail inequalities for random tries and Patricia trees,” *Journal of Computational and Applied Mathematics*, vol. 142, pp. 27–37, 2002.
- B. Efron and C. Stein, “The jackknife estimate of variance,” *Annals of Statistics*, vol. 9, pp. 586–596, 1981.
- P. Flajolet and J. M. Steyaert, “A branching process arising in dynamic hashing, trie searching and polynomial factorization,” in: *Lecture Notes in Computer Science*, vol. 140, pp. 239–251, Springer-Verlag, New York, 1982.
- P. Flajolet, “On the performance evaluation of extendible hashing and trie search,” *Acta Informatica*, vol. 20, pp. 345–369, 1983.
- P. Flajolet, M. Régnier, and D. Sotteau, “Algebraic methods for trie statistics,” *Annals of Discrete Mathematics*, vol. 25, pp. 145–188, 1985.
- P. Flajolet and C. Puech, “Tree structure for partial match retrieval,” *Journal of the ACM*, vol. 33, pp. 371–407, 1986.
- P. Flajolet and R. Sedgewick, “Digital search trees revisited,” *Siam Journal on Computing*, vol. 15, pp. 748–767, 1986.
- P. Flajolet and B. Vallée, “Continued fraction algorithms, functional operators, and structure constants,” *Theoretical Computer Science*, vol. 194, pp. 1–34, 1998.
- E. H. Fredkin, “Trie memory,” *Communications of the ACM*, vol. 3, pp. 490–500, 1960.
- G. H. Gonnet, “Expected length of the longest probe sequence in hash code searching,” *Journal of the ACM*, vol. 28, pp. 289–304, 1981.

- W. Hoeffding, “Probability inequalities for sums of bounded random variables,” *Journal of the American Statistical Association*, vol. 58, pp. 13–30, 1963.
- P. Jacquet and M. Régnier, “Trie partitioning process: limiting distributions,” in: *CAAP 86*, (edited by P. Franchi-Zannettacci), vol. 214, pp. 196–210, Lecture Notes in Computer Science, Springer-Verlag, Berlin, 1986.
- P. Jacquet and M. Régnier, “Normal limiting distribution for the size and the external path length of tries,” Technical Report, INRIA, Rocquencourt, France, 1989.
- P. Jacquet, W. Szpankowski, and Analysis of digital tries with Markovian dependency, *IEEE Transactions on Information Theory*, vol. IT-37, pp. 1470–1475, 1991.
- P. Kirschenhofer and H. Prodinger, “Some further results on digital trees,” in: *Lecture Notes in Computer Science*, vol. 226, pp. 177–185, Springer-Verlag, Berlin, 1986.
- P. Kirschenhofer, H. Prodinger, and W. Szpankowski, “On the balance property of PATRICIA tries: external path length viewpoint,” *Theoretical Computer Science*, vol. 68, pp. 1–17, 1989.
- P. Kirschenhofer, H. Prodinger, and W. Szpankowski, “On the variance of the external path length in a symmetric digital trie,” *Discrete Applied Mathematics*, vol. 25, pp. 129–143, 1989.
- P. Kirschenhofer, H. Prodinger, and W. Szpankowski, “Multidimensional digital searching and some new parameters in tries,” *International Journal of Foundations of Computer Science*, vol. 4, pp. 69–84, 1993.
- D. E. Knuth, *The Art of Computer Programming, Vol. 3 : Sorting and Searching*, Addison-Wesley, Reading, Mass., 1973.
- M. Ledoux, “On Talagrand’s deviation inequalities for product measures,” *ESAIM: Probability and Statistics*, vol. 1, pp. 63–87, 1996a.
- M. Ledoux, “Isoperimetry and gaussian analysis,” in: *Lectures on Probability Theory and Statistics*, (edited by P. Bernard), pp. 165–294, Ecole d’Eté de Probabilités de St-Flour XXIV-1994, 1996b.
- C. McDiarmid, “On the method of bounded differences,” in: *Surveys in Combinatorics*, (edited by J. Siemons), vol. 141, pp. 148–188, London Mathematical Society Lecture Note Series, Cambridge University Press, 1989.
- C. McDiarmid, “Concentration,” in: *Probabilistic Methods for Algorithmic Discrete Mathematics*, (edited by M. Habib and C. McDiarmid and J. Ramirez-Alfonsin and B. Reed), pp. 195–248, Springer, New York, 1998.
- H. Mendelson, “Analysis of extendible hashing,” *IEEE Transactions on Software Engineering*, vol. 8, pp. 611–619, 1982.
- D. R. Morrison, “PATRICIA — Practical Algorithm To Retrieve Information Coded in Alphanumeric,” *Journal of the ACM*, vol. 15, pp. 514–534, 1968.

- J. A. Orenstein, “Multidimensional tries used for associative searching,” Technical Report, School of Computer Science, McGill University, Montreal, 1982.
- B. Pittel, “Asymptotical growth of a class of random trees,” *Annals of Probability*, vol. 13, pp. 414–427, 1985.
- B. Pittel, “Path in a random digital tree: limiting distributions,” *Advances in Applied Probability*, vol. 18, pp. 139–155, 1986.
- B. Pittel and H. Rubin, “How many random questions are necessary to identify  $n$  distinct objects?,” *Journal of Combinatorial Theory*, vol. A55, pp. 292–312, 1990.
- B. Pittel, “On the height of PATRICIA search tree,” ORSA/TIMS Special Interest Conference on Applied Probability in the Engineering, Information and Natural Sciences, Monterey, CA, 1991.
- C. Puech and H. Yahia, “Quadrees, octrees, hyperoctrees: a unified analytical approach to tree data structures used in graphics, geometric modeling and image processing,” in: *Proceedings of the Symposium on Computational Geometry*, pp. 272–280, ACM, New York, 1985.
- B. Rais, P. Jacquet, and W. Szpankowski, “A limiting distribution for the depth in Patricia tries,” *SIAM Journal on Discrete Mathematics*, vol. 6, pp. 197–213, 1993.
- M. Régnier, “On the average height of trees in digital searching and dynamic hashing,” *Information Processing Letters*, vol. 13, pp. 64–66, 1981.
- M. Régnier, “Trie hashing analysis,” in: *Proceedings of the Fourth International Conference on Data Engineering*, pp. 377–381, IEEE, Los Angeles, 1988.
- M. Régnier and P. Jacquet, “New results on the size of tries,” *IEEE Transactions on Information Theory*, vol. IT-35, pp. 203–205, 1989.
- W. Schachinger, “The variance of a partial match retrieval in a multidimensional symmetric trie,” *Random Structures and Algorithms*, vol. 7, pp. 81–95, 1995.
- W. Schachinger, “Limiting distributions for the costs of partial match retrievals in multidimensional tries,” *Random Structures and Algorithms*, vol. 17, pp. 428–459, 2000.
- W. Schachinger, “Concentration of distribution results for trie-based sorting of continued fractions,” Technical Report, Universität Wien, Austria, 2002.
- W. Szpankowski, “Some results on  $V$ -ary asymmetric tries,” *Journal of Algorithms*, vol. 9, pp. 224–244, 1988.
- W. Szpankowski, “Digital data structures and order statistics,” in: *Algorithms and Data Structures: Workshop WADS ’89 Ottawa*, vol. 382, pp. 206–217, Lecture Notes in Computer Science, Springer-Verlag, Berlin, 1989.
- W. Szpankowski, “Patricia trees again revisited,” *Journal of the ACM*, vol. 37, pp. 691–711, 1990.

- W. Szpankowski, "On the height of digital trees and related problems," *Algorithmica*, vol. 6, pp. 256–277, 1991.
- W. Szpankowski and C. Knessl, "Heights in generalized tries and PATRICIA tries," in: *LATIN'2000*, pp. 298–307, Lecture Notes in Computer Science 1776, 2000.
- M. Talagrand, "An isoperimetric theorem on the cube and the Kintchine-Kahane inequalities," *Proceedings of the American Mathematical Society*, vol. 104, pp. 905–909, 1988.
- M. Talagrand, "Isoperimetry and integrability of the sum of independent Banach-space valued random variables," *Annals of Probability*, vol. 17, pp. 1546–1570, 1989.
- M. Talagrand, "Sample unboundedness of stochastic processes under increment conditions," *Annals of Probability*, vol. 18, pp. 1–49, 1990.
- M. Talagrand, "A new isoperimetric inequality and the concentration of measure phenomenon," in: *Geometric Aspects of Functional Analysis*, (edited by J. Lindenstrauss and V. D. Milman), vol. 146, pp. 4–137, Lecture Notes in Mathematics, Springer-Verlag, 1991a.
- M. Talagrand, "A new isoperimetric inequality for product measure and the tails of sums of independent random variables," *Geometric and Functional Analysis*, vol. 1, pp. 211–223, 1991b.
- M. Talagrand, "Isoperimetry, logarithmic Sobolev inequalities on the discrete cube, and Margulis' graph connectivity theorem," *Geometric and Functional Analysis*, vol. 3, pp. 295–314, 1993a.
- M. Talagrand, "New gaussian estimates for enlarged balls," *Geometric and Functional Analysis*, vol. 3, pp. 502–526, 1993b.
- M. Talagrand, "Sharper bounds for gaussian and empirical processes," *Annals of Probability*, vol. 22, pp. 28–76, 1994.
- M. Talagrand, "Concentration of measure and isoperimetric inequalities in product spaces," *Inst. Hautes Etudes Sci. Publ. Math.*, vol. 81, pp. 73–205, 1995.
- M. Talagrand, "New concentration inequalities in product spaces," *Inventiones Mathematicae*, vol. 126, pp. 505–563, 1996a.
- M. Talagrand, "A new look at independence," *Annals of Probability*, vol. 24, pp. 1–34, 1996b.
- A. W. van der Vaart and J. A. Wellner, *Weak Convergence and Empirical Processes*, Springer-Verlag, New York, 1996.