

ON THE NUMBER OF DESCENDANTS AND ASCENDANTS IN RANDOM SEARCH TREES*

CONRADO MARTÍNEZ

Departament de Llenguatges i Sistemes Informàtics,
Polytechnical University of Catalonia,
Pau Gargallo 5, E-08028 Barcelona, Spain.
email: Conrado.Martinez@lsi.upc.es
www: <http://www-lsi.upc.es/~conrado/home.html>

ALOIS PANHOLZER

Institut für Algebra und Diskrete Mathematik,
Technical University of Vienna,
Wiedner Hauptstrasse 8–10,
A-1040 Vienna, Austria.
email: e9125354@fbma.tuwien.ac.at

HELMUT PRODINGER

Institut für Algebra und Diskrete Mathematik,
Technical University of Vienna,
Wiedner Hauptstrasse 8–10,
A-1040 Vienna, Austria.
email: Helmut.Prodinger@tuwien.ac.at
www: <http://info.tuwien.ac.at/theoinf/proding.htm>

Submitted: January 7, 1997; Accepted: March 26, 1998.

ABSTRACT. The number of descendants of a node in a binary search tree (BST) is the size of the subtree having this node as a root; the number of ascendants is the number of nodes on the path connecting this node with the root. Using a purely combinatorial approach (generating functions and differential equations) we are able to extend previous results. For the number of descendants we get explicit formulae for all moments; for the number of ascendants, which is harder, we get the variance.

A natural extension of binary search trees occurs when performing local reorganisations. Poblete and Munro have already analyzed some aspects of these locally balanced binary search trees (LBSTs). Here, we relate these structures with the performance of median-of-three Quicksort. We get as new results the variances for ascendants and descendants in this setting.

If the rank of the node itself is picked at random (“grand averages”), the corresponding parameters only depend on the size n . In this instance, we get all the moments for the descendants (BST and LBST), as well as the probabilities. For ascendants (LBST), we get the variance and (in principle) the higher moments, as well as the (normal) limiting distribution.

The emphasis is on explicit formulae, and these are sometimes quite involved. Thus, in some instances, we have decided to state abridged versions in the paper and collect the long forms into an appendix that can be downloaded from the URLs http://info.tuwien.ac.at/theoinf/abstract/abs_120.htm and <http://www.lsi.upc.es/~conrado/research/>.

AMS Subject Classification. 05A15 (primary) 05C05, 68P10 (secondary)

* This research was partly done while the third author was visiting the CRM (Centre de Recerca Matemàtica, Institut d’Estudis Catalans). The first author was supported by the ESPRIT Long Term Research Project ALCOM IT (contract no. 20244). The second author was supported by the FWF Project 12599-MAT. All 3 authors are supported by the Project 16/98 of Acciones Integradas 1998/99.

The appendix of this paper with all the outsize expressions is downloadable from the URLs http://info.tuwien.ac.at/theoinf/abstract/abs_120.htm and <http://www.lsi.upc.es/~conrado/research/>.

1. INTRODUCTION

Binary search trees are among the most important and commonly used data structures, their applications spanning a wide range of the areas of Computer Science. Standard binary search trees (BSTs, for short) are still the subject of active research, see for instance the recent articles [2, 28]. Deepening our knowledge about binary search trees is interesting in its own; moreover, most of this knowledge can be translated and applied to other data structures such as heap ordered trees, k -d-trees [33], and to important algorithms like quicksort and Hoare's Find algorithm for selection (also known as quickselect) [12, 13, 30, 31].

We assume that the reader is already familiar with binary search trees and the basic algorithms to manipulate them [20, 31, 9]. Height and weight-balanced versions of the binary search trees, like AVL and red-black trees [1, 11], have been proposed and find many useful applications, since all of them guarantee good worst-case performance of both searches and updates.

Locally balanced search trees (LBSTs) were introduced by Bell [4] and Walker and Wood [34], and thoroughly analyzed by Poblete and Munro in [27]. LBSTs have been proposed as an alternative to more complex balancing schemes for search trees. In these search trees, only local rebalancing is made; after each insertion, local rebalancing is applied to ensure that all subtrees of size 3 in the tree are complete¹. The basic idea of the heuristic is that the construction of poorly balanced trees becomes less likely. A similar idea, namely, selecting a sample of 3 elements and taking the median of the sample as the pivot element for partitioning in algorithms like quicksort and quickselect has been shown to yield significant improvements in theory and practice [30, 17].

Random search trees, either random BSTs or random LBSTs, are search trees built by performing n random insertions into an initially empty tree [20, 24]. An insertion of a new element into a search tree of size k is said to be random, if the new element falls with equal probability into any of the $k + 1$ intervals defined by the k keys already present in the tree (equivalently, the new element replaces any of the $k + 1$ external nodes in the tree with equal probability). Random search trees can also be defined as the result of the insertion of the elements of a random permutation of $\{1, \dots, n\}$ into an initially empty tree.

Ascendants and descendants of the j^{th} internal node of a random search tree of size n are denoted $A_{n,j}$ and $D_{n,j}$, respectively. Besides the two aforementioned random variables, we also consider other random variables: the number of descendants D_n and the number of ascendants A_n of a randomly chosen internal node in a random search tree of size n . This corresponds to averaging $D_{n,j}$ and $A_{n,j}$ over j . We remark, that all the distributions, as well as the expectations $\mathbb{E}[X]$ and probabilities $\mathbb{P}[X]$ are induced by the creation process of the random search trees (BSTs resp. LBSTs). The number of descendants and the number of ascendants in random BSTs have been investigated in several previous works ([3, 5, 23, 22, 21]). The number of ascendants of a random node in a random LBST has been studied in [27, 26].

We define the number of descendants $D_{n,j}$ as the size of the subtree rooted at the j^{th} node, so we count the j^{th} node as a descendant of itself. The number of ascendants $A_{n,j}$ is the number of internal nodes in the path from the root of the tree to the j^{th} node, both included. It is worth mentioning the following symmetry property (which is very easy to prove) for the random variables we are going to consider.²

¹The generalization of the local rebalancing heuristic to subtree sizes larger than 3 is straightforward.

²We remark, that here and in the sequel equalities between random variables are equalities in distribution, which is often denoted by $\stackrel{d}{=}$.

Proposition 1.1. For any $n > 0$ and any $1 \leq j \leq n$,

$$\begin{aligned} D_{n,j} &= D_{n,n+1-j}, \\ A_{n,j} &= A_{n,n+1-j}. \end{aligned}$$

The performance of a successful search is obviously proportional to the number of ascendants of the sought internal node. The next proposition states this relation, as well as other interesting relationships that hold for both random BSTs and random LBSTs.

Proposition 1.2. Consider a random search tree of size n and let

$$\begin{aligned} S_{n,j} &= \# \text{ of comparisons in a successful search for the } j^{\text{th}} \text{ element,} \\ S_n &= \# \text{ of comparisons in a successful search for a randomly chosen element,} \\ U_n &= \# \text{ of comparisons in an unsuccessful search for a randomly chosen external node,} \\ P_{n,j} &= \text{depth of the } j^{\text{th}} \text{ element,} \\ I_n &= \sum_{1 \leq j \leq n} P_{n,j} = \text{internal path length,} \end{aligned}$$

Then,

$$\begin{aligned} S_{n,j} &= P_{n,j} + 1 = A_{n,j}, \\ S_n &= A_n, \\ \mathbb{E}[U_n] &= \frac{n}{n+1} (1 + \mathbb{E}[A_n]), \\ \mathbb{E}[I_n] &= n (\mathbb{E}[A_n] - 1), \\ \mathbb{E}[A_n] &= \mathbb{E}[D_n]. \end{aligned}$$

There is also a close relationship between the performance of quickselect [12, 19, 17] and the number of ascendants.

Proposition 1.3. Let $F_{n,j}$ be the number of recursive calls made by quickselect to select the j^{th} element out of n elements. Then

$$F_{n,j} = A_{n,j}.$$

If we consider $A_{n,j}$ in random BSTs, then this corresponds to the selection of the pivots at random in each phase of quickselect. If we consider $A_{n,j}$ in random LBSTs, then the proposition applies for the variant of quickselect that uses the median of a random sample of three elements as the pivot in each partitioning phase.

The study of the number of descendants has applications in the context of paged trees (see for instance [20, 14]). A paged binary search tree with *page capacity* b stores all its subtrees of size $\leq b$ (possibly empty) in pages; typically, the pages reside in secondary memory and the elements within a page are not organized as search trees (see Figure 1: the pagination of the search tree at the left is indicated using dashed lines; a more “realistic” representation of the same tree appears at its right).

Let $\mathcal{P}_n^{(b)}$ be the number of pages in a random search tree of size n with page capacity b . It is obvious that $\mathcal{P}_n^{(b)} = \mathcal{I}_n^{(b)} + 1$, where $\mathcal{I}_n^{(b)}$ is the number of internal nodes that are the root of a subtree that contains more than b items. In other words, in a paged search tree, we have external nodes (*pages*) that may contain up to b keys; if $\mathcal{P}_n^{(b)}$ is the number of external nodes or pages in a paged search tree, then $\mathcal{I}_n^{(b)} = \mathcal{P}_n^{(b)} - 1$ is the number of internal nodes in the tree, and these internal nodes are in one-to-one correspondance with the internal nodes with $> b$ descendants in the non-paged search tree.

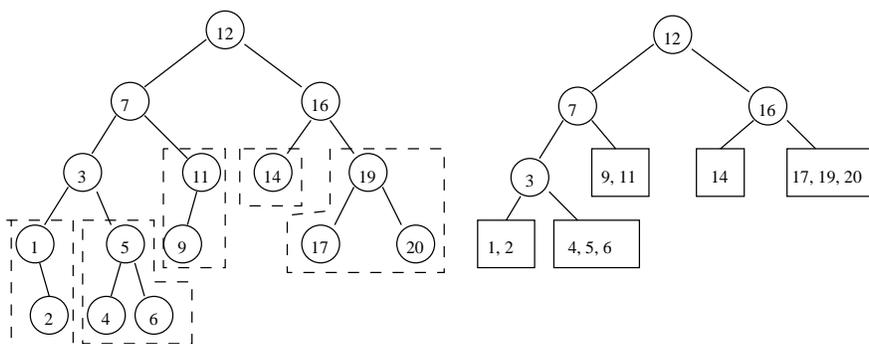


FIGURE 1. A paged binary search tree with page capacity $b = 3$

Proposition 1.4. For all n , and for any constant $b \geq 1$,

$$\mathbb{E} \left[\mathcal{P}_n^{(b)} \right] = n \mathbb{P} [D_n > b] + 1.$$

Proof. Let δ_j be the indicator random variable for the predicate “the j^{th} element has more than b descendants.”. Then $\mathcal{I}_n^{(b)} = \sum_{1 \leq j \leq n} \delta_j$. The proposition follows taking expectations in both sides of this equation, because of the linearity of expectations and $\mathbb{E} [\delta_j] = \mathbb{P} [D_{n,j} > b]$. \square

Results about the probabilistic behavior of the number of descendants are also useful in the analysis of the performance of quicksort if recursive calls are not made on small subfiles (say, of size $\leq b$).

Proposition 1.5. Let $C_n^{(b)}$ and $R_n^{(b)}$ be the number of comparisons³ and the number of partitions made by quicksort to sort n elements, when the recursion halts on subfiles of size $\leq b$. Notice that standard quicksort corresponds to the case where $b = 1$. Then

$$\begin{aligned} \mathbb{E} \left[R_n^{(b)} \right] &= n \mathbb{P} [D_n > b], \\ \mathbb{E} \left[C_n^{(b)} \right] &= n (\mathbb{E} [D_n] - 1) - n \sum_{1 \leq m \leq b} (m - 1) \mathbb{P} [D_n = m]. \end{aligned}$$

The strategy for the selection of pivots is related with the type of random search trees that we consider: for BSTs, we have selection of pivots at random; for LBSTs, we have that the pivots are the medians of random samples of three elements.

Proof. It is well known that we can associate to each particular execution of quicksort a binary search tree: the root contains the pivot element of the first stage, and the left and right subtrees are recursively built for the elements smaller and larger than the pivot, respectively. Each internal node in the search tree corresponds to a recursive call to quicksort. We will make a partitioning of a given subfile if and only if the subfile contains $> b$ elements, i.e. the corresponding internal node has $> b$ descendants, and the claim in the proposition follows.

On the other hand, let ϵ_j be the number of comparisons made between the j^{th} element and other elements, during the partition where the j^{th} element was selected as a pivot. Clearly, if $D_{n,j} \leq b$ then $\epsilon_j = 0$, since no recursive call will be made that chooses the j^{th} element as a pivot. On the other hand, if $D_{n,j} > b$, the j^{th} element will be compared with each of its descendants (except itself) in the associated search tree. Hence, $\mathbb{E} [\epsilon_j] = \sum_{m=b+1}^n (m - 1) \mathbb{P} [D_{n,j} = m]$. We need only to sum over j to get the desired result. \square

³We only count those made during the partitioning phases.

	BST		LBST	
	Of a given node	Of a random node	Of a given node	Of a random node
Ascendants	Average [3], variance*	Probability, moments, limit distribution [23, 5, 22, 18]	Average [17], variance*	Average, variance [27]*, higher order moments, PGF, limit distribution*
Descendants	Probability, moments [21]*	Probability, moments [21]*	PGF, average, variance*	Probability, moments*

TABLE 1. Summary of previous works and the results of this paper.

The structure of the paper is as follows. We start with an overview of some basic facts about generating functions and, in particular, about probability generating functions (Section 2).

In Section 3 we develop the main steps of our approach, taking the analysis of the number of descendants in random BSTs as a first introductory example. We provide here alternative derivations to the results of Lent [21], finding the probability that the j^{th} node in a random BST of size n has m descendants (Theorem 3.1). We also find exact and asymptotic values for all ordinary moments, including the expected value and variance (Theorem 3.2). Then we analyze the number of descendants of a random node, obtaining the probability that $D_n = m$, as well as the moments of D_n (Theorems 3.3 and 3.2).

The remaining sections are devoted to the analysis of the number of ascendants and descendants in random LBSTs. In Section 5 we formally define LBSTs and give an equivalent characterization of the model of randomness which is more suitable to our purposes.

Among our new results, in Section 6 we derive an explicit form for the generating function of the probability distribution of $D_{n,j}$ (Theorem 6.1) and closed formulæ for the average (Theorem 6.2) and the second factorial moment (Theorem 6.3). Moreover, we find the probability distribution of D_n (Theorem 6.4) and all its moments (Theorem 6.5).

In Section 7, we compute $\mathbb{E}[A_{n,j}]$, the average number of ascendants of the j^{th} node in a random LBST of size n (Theorem 7.1). We are also able to compute the PGF of A_n , the number of ascendants of a random node (Theorem 7.2), as well as all its moments (Theorems 7.4 and 7.5), thus extending the results of Poblete and Munro [27].

The results of previous works and the new results in this paper are summarized in Table 1. Entries corresponding to new results in this paper and to alternative derivations of previous results are marked by ‘*’.

2. MATHEMATICAL PRELIMINARIES

We start recalling the definition of *generating function*, for the reader’s convenience. Given a sequence $\{a_n\}_{n \geq 0}$ its generating function $A(z)$ is the formal power series

$$A(z) = \sum_{n \geq 0} a_n z^n.$$

As usual, $[z^n]A(z)$ denotes the coefficient of z^n in $A(z)$ (the n^{th} coefficient of $A(z)$). Excellent sources of information about generating functions and their applications to combinatorics and the analysis of algorithms are [35, 33, 32, 20].

We make extensive use in this paper of probability generating functions (PGFs) as well as multivariate generating functions whose coefficients are PGFs themselves. We define them in turn. Given a discrete random variable X , its probability generating function $X(z)$ is

$$X(z) = \sum_m \mathbb{P}[X = m] z^m.$$

If we assume further that $X \geq 0$ and let $p_m = \mathbb{P}[X = m]$, the PGF of the random variable X is nothing but the ordinary generating function of the sequence $\{p_m\}_{m \geq 0}$. We list now a few important, although elementary, properties of PGFs.

Proposition 2.1. *For any discrete random variable X , its probability generating function $X(z)$ satisfies:*

1. $X(1) = 1$.
2. $X'(1) = \left. \frac{dX}{dz} \right|_{z=1} = \mathbb{E}[X]$.
3. $X^{(s)}(1) = \left. \frac{d^s X}{dz^s} \right|_{z=1} = \mathbb{E}[X^{\underline{s}}]$, where $X^{\underline{s}}$ denotes the s^{th} falling factorial of X , that is, $X^{\underline{s}} = X(X-1)\dots(X-s+1)$. The quantity $\mathbb{E}[X^{\underline{s}}]$ is customarily called the s^{th} factorial moment of the random variable X . Ordinary and central moments may be recovered from factorial moments quite easily. For instance, if $\mu = \mathbb{E}[X]$, the variance of X is given by

$$\mathbb{V}[X] = \mathbb{E}[(X - \mu)^2] = \mathbb{E}[X^2] + \mathbb{E}[X] - \mathbb{E}[X]^2.$$

Since we will mostly deal with families of random variables, with two (n and j) or one (n) index, we will systematically work with multivariate generating functions of these families. For instance, if we were interested in the family $\{X_{n,j}\}_{1 \leq j \leq n}$, we would introduce a generating function $X(z, u, v)$ in three variables, such that the coefficient of $z^n u^j v^m$ in $X(z, u, v)$ is the probability that $X_{n,j}$ is m . Thus

$$X(z, u, v) = \sum_{n,j,m} \mathbb{P}[X_{n,j} = m] z^n u^j v^m, \tag{1}$$

where the indices of summation n, j and m run in the appropriate ranges (or we assume that $\mathbb{P}[X_{n,j} = m]$ is 0 whenever $n < 1, j < 1, j > n$ or $m < 0$). Notice that, by definition, $[z^n u^j]X(z, u, v)$ is the PGF of the random variable $X_{n,j}$, and $[z^n u^j v^m]X(z, u, v) = \mathbb{P}[X_{n,j} = m]$.

For technical reasons that will be clearer later, we will also use sometimes the derivative w.r.t. z of such a multivariate generating function. We will introduce then

$$\begin{aligned} X_z(z, u, v) &= \frac{\partial}{\partial z} \sum_{n,j,m} \mathbb{P}[X_{n,j} = m] z^n u^j v^m \\ &= \sum_{n,j,m} n \mathbb{P}[X_{n,j} = m] z^{n-1} u^j v^m \end{aligned}$$

rather than the more natural definition given in Equation (1). This means that once we were able to extract coefficients from such a generating function, let us say the coefficient of $z^{n-1} u^j v^m$, we must divide by n to obtain $\mathbb{P}[X_{n,j} = m]$.

Furthermore, we are also interested in investigating all the moments of the random variables: mean, variance, and higher order moments. We differentiate the generating function $X(z, u, v)$ s times with respect to v and let $v = 1$, to get the generating function for the s^{th} factorial moments, i.e.

$$\mathcal{X}^{(s)}(z, u) = \left. \frac{\partial^s X(z, u, v)}{\partial v^s} \right|_{v=1}, \quad s \geq 1. \tag{2}$$

Recall that $[z^n u^j] \mathcal{X}^{(s)}(z, u) = \mathbb{E}[X_{n,j}^{\underline{s}}]$.

Grand averages correspond to the situation where the *rank* —the parameter j in $X_{n,j}$ — is random itself. More precisely, let $X_n \equiv X_{n,Z_n}$, where Z_n is a uniformly distributed random variable

in $\{1, \dots, n\}$. Then X_n is the grand average of the random variables $X_{n,1}, \dots, X_{n,n}$. It follows that

$$\mathbb{P}[X_n = m] = \frac{1}{n} \sum_{1 \leq j \leq n} \mathbb{P}[X_{n,j} = m]. \quad (3)$$

We remark that $X_n \neq \frac{1}{n}(X_{n,1} + \dots + X_{n,n})$, even if the $X_{n,j}$'s are independent.

Unless we are dealing with a differentiated version of the generating function $\mathsf{X}(z, u, v)$, we have

$$\overline{\mathsf{X}}(z, v) = \mathsf{X}(z, 1, v) = \sum_{n,m} z^n v^m \sum_{1 \leq j \leq n} \mathbb{P}[X_{n,j} = m]. \quad (4)$$

Thus the coefficient $[z^n v^m] \overline{\mathsf{X}}(z, v)$, divided by n , is the probability that X_n is m . In the case that $\mathsf{X}_z(z, u, v)$ were a differentiated generating function, then we should divide the coefficient $[z^{n-1} v^m] \overline{\mathsf{X}}_z(z, v)$ by n^2 . Finally, computing the derivatives of $\mathsf{X}(z, v)$ w.r.t. v and setting $v = 1$ yields the generating functions for the factorial moments of the grand average X_n .

The main steps of the systematic procedure that we will follow are thus:

1. Set up a recurrence for $\mathbb{P}[X_{n,j} = m]$;
2. Translate the recurrence to a functional equation over the corresponding generating function $\mathsf{X}(z, u, v)$;
3. Solve the functional equation;
4. Extract the coefficients of $\mathsf{X}(z, u, v)$;
5. Repeatedly differentiate $\mathsf{X}(z, u, v)$ w.r.t. v and set $v = 1$; extract the coefficients to get the factorial moments of $X_{n,j}$;
6. Set $\overline{\mathsf{X}}(z, v) = \mathsf{X}(z, 1, v)$ and repeat steps 4 and 5 for $\overline{\mathsf{X}}(z, v)$.

In practice, the procedure might fail for several reasons. Typically, because we are not able to solve the equation at step 3 or to extract the coefficients of a given generating function. Although we have (almost) not used them in this paper, the reader should be aware of the existing powerful techniques to extract asymptotic information about the coefficients of a generating function if we know its behaviour near its singularities or in some case, even if we only know the functional equation satisfied by the generating function [33, 6]. Also, if we are not able to solve and get an explicit form for $\mathsf{X}(z, u, v)$, we can still differentiate w.r.t. to v or set $u = 1$ and try to solve the (easier) resulting differential equations, to get information about the moments or the grand average.

The functional equations that arise in our study are linear partial differential equations of the first (BSTs) and of the second (LBSTs) order. The former can be solved, in principle, by quadrature through the variation of constant —actually, functions in u and v — method. For the second order differential equations, the theory of hypergeometric differential equations comes into play [16].

Nowadays, most of the necessary mathematical knowledge is embodied into modern computer algebra systems. In our case, MAPLE needed little or no assistance to solve the differential equations that we had.

The last step, that of extracting coefficients in exact form, was, at large, the least systematic and mechanical one. A great deal of combinatorial identities, inspired guessing and patience was needed. Standard MAPLE tools like the function `interp` or the `GFUN` package [29] proved also to be useful. However,

once the solution is obtained, it is just a matter of minutes to check its correctness. It is quite difficult to provide a detailed and ordered description of the methods that we used to extract coefficients from generating functions. As a result, the paper contains only some hints here and there, while some claims are just stated without further explanation.

3. THE NUMBER OF DESCENDANTS IN RANDOM BSTs

The number of the descendants $D_{n,j}$ of the j^{th} node of a BST of size n is recursively computed as the number of descendants in the left subtree of the j^{th} node, plus the number of descendants in its right subtree, plus one (to count the j^{th} node itself). The probability that $D_{n,j} = m$ is computed conditioning on the events “the rank of the root is k ,” that means the root is the k^{th} node of a search tree. Recall that, for a random BST of size n , the rank of the root is k with probability $1/n$, for $k = 1, \dots, n$. Using the recursive definition of $D_{n,j}$ we have

$$\begin{aligned} \mathbb{P}[D_{n,j} = m] &= \sum_{k=1}^n \mathbb{P} \left[D_{n,j} = m \mid \text{the root is the } k^{\text{th}} \text{ element} \right] \\ &\quad \times \mathbb{P} \left[\text{the root is the } k^{\text{th}} \text{ element} \right] \\ &= \frac{1}{n} \llbracket m = n \rrbracket + \frac{1}{n} \sum_{k=1}^{j-1} \mathbb{P} [D_{n-k,j-k} = m] + \frac{1}{n} \sum_{k=j+1}^n \mathbb{P} [D_{k-1,j} = m], \end{aligned} \tag{5}$$

where $\llbracket P \rrbracket$ is 1 if P is true and 0 otherwise [10].

This recursion translates nicely into a functional equation over the generating function for the family of random variables $\{D_{n,j}\}$. Solving the functional equation and extracting coefficients of the generating function, we get the following theorem, which was already found by Lent [21] using probabilistic techniques.

Theorem 3.1. *The probability that the j^{th} internal node of a random binary search tree of size n has m descendants is, assuming that $j \leq n + 1 - j$,*

$$\mathbb{P}[D_{n,j} = m] = \begin{cases} \frac{2}{(m+1)(m+2)} & \text{for } 1 \leq m < j, \\ \frac{1}{(m+1)(m+2)} \left(1 + \frac{2j}{m}\right) & \text{for } j \leq m < n + 1 - j, \\ \frac{2(n+1)}{m(m+1)(m+2)} & \text{for } n + 1 - j \leq m < n, \\ \frac{1}{n} & \text{for } m = n. \end{cases}$$

For the cases where $j > n + 1 - j$ we can use the symmetry on j and $n + 1 - j$ (Proposition 1.1) to compute the corresponding probabilities.

Also, the distribution function for $D_{n,j}$ is

$$\mathbb{P}[D_{n,j} \leq m] = \begin{cases} \frac{m}{m+2} & \text{for } 1 \leq m < j, \\ \frac{m+1}{m+2} - \frac{j}{(m+1)(m+2)} & \text{for } j \leq m < n + 1 - j, \\ \frac{m^2 + 3m + 1 - n}{(m+1)(m+2)} & \text{for } n + 1 - j \leq m < n, \\ 1 & \text{for } m = n. \end{cases}$$

Proof. We start defining the generating function

$$D(z, u, v) = \sum_{1 \leq j, m \leq n} \mathbb{P}[D_{n,j} = m] z^n u^j v^m.$$

Multiplying both sides of (5) by $nz^{n-1}u^jv^m$ and summing for all $n \geq 1$, $1 \leq j \leq n$ and $m \geq 1$, yields

$$\frac{\partial D}{\partial z} = \frac{uD}{1-uz} + \frac{D}{1-z} + \frac{uv}{(1-vz)(1-uvz)},$$

$$D(0, u, v) = 0. \tag{6}$$

The solution to the differential equation above is relatively simple

$$D(z, u, v) = \frac{uz}{v(1-z)(1-uz)}$$

$$- \frac{u(1-v)(v-u)}{(1-z)(1-uz)v^2(1-u)} \log \frac{1}{1-vz}$$

$$- \frac{(1-v)(1-uv)}{(1-z)(1-uz)v^2(1-u)} \log \frac{1}{1-uvz}. \tag{7}$$

The statement of the theorem follows after extracting the coefficient $[z^n u^j v^m]D(z, u, v)$. □

The explicit and simple form of the trivariate generating function in Theorem 3.1 allows us to compute all the moments *explicitly*. It is convenient to deal with a sort of shifted factorial moments; the ordinary moments can be computed by linear combinations of the shifted factorial ones.

Theorem 3.2. *Let $d_{n,j}^{(s)} = \mathbb{E}[(D_{n,j} + 2)^s]$ and $d_{n,j} = d_{n,j}^{(1)}$, where $D_{n,j}$ denotes the number of descendants of the j^{th} internal node in a random binary search tree of size n . For all $n > 0$ and all $1 \leq j \leq n$,*

1. $d_{n,j} = H_j + H_{n+1-j} + 1$,
2. $d_{n,j}^{(2)} = 2(n+1)H_n - 2jH_j - 2(n+1-j)H_{n+1-j} + 2(n+2)$.
3. For all $s \geq 3$,

$$d_{n,j}^{(s)} = \frac{s}{s-2}(n+1)^{s-1} - \frac{s}{(s-1)(s-2)} \left[j^{s-1} + (n+1-j)^{s-1} \right].$$

Proof. We begin by introducing

$$\mathcal{D}^{(s)}(z, u) = \left. \frac{\partial^s (v^2 D(z, u, v))}{\partial v^s} \right|_{v=1},$$

and hence its coefficients are

$$d_{n,j}^{(s)} = [z^n u^j] \mathcal{D}^{(s)}(z, u) = \mathbb{E}[(D_{n,j} + 2)^s].$$

The shifted moments are particularly easy to obtain, since the coefficients of $\mathcal{D}^{(s)}(z, u)$ that we seek are linear combinations of the coefficients of the next generating functions:

$$\frac{\partial^s}{\partial v^s} \log \frac{1}{1-vz} \Big|_{v=1} = (s-1)! \left(\frac{z}{1-z} \right)^s,$$

$$\frac{\partial^s}{\partial v^s} v \log \frac{1}{1-vz} \Big|_{v=1} = (s-1)! \left(\frac{z}{1-z} \right)^s + s(s-2)! \left(\frac{z}{1-z} \right)^{s-1},$$

$$\frac{\partial^s}{\partial v^s} v^2 \log \frac{1}{1-vz} \Big|_{v=1} = (s-1)! \left(\frac{z}{1-z} \right)^s + 2s(s-2)! \left(\frac{z}{1-z} \right)^{s-1}$$

$$+ s(s-1)(s-3)! \left(\frac{z}{1-z} \right)^{s-2},$$

$$\frac{\partial^s}{\partial v^s} \log \frac{1}{1-uvz} \Big|_{v=1} = (s-1)! \left(\frac{uz}{1-uz} \right)^s,$$

$$\begin{aligned} \frac{\partial^s}{\partial v^s} v \log \frac{1}{1-uvz} \Big|_{v=1} &= (s-1)! \left(\frac{uz}{1-uz}\right)^s + s(s-2)! \left(\frac{uz}{1-uz}\right)^{s-1}, \\ \frac{\partial^s}{\partial v^s} v^2 \log \frac{1}{1-uvz} \Big|_{v=1} &= (s-1)! \left(\frac{uz}{1-uz}\right)^s + 2s(s-2)! \left(\frac{uz}{1-uz}\right)^{s-1} \\ &\quad + s(s-1)(s-3)! \left(\frac{uz}{1-uz}\right)^{s-2}. \end{aligned}$$

We might additionally observe that for all $n \geq 0$ and $1 \leq j \leq n$

$$\begin{aligned} [z^n u^j] \frac{1}{(1-z)^{s+1}(1-uz)(1-u)} &= \binom{s+n+1}{s+1} - \binom{s+n-j}{s+1}, \\ [z^n u^j] \frac{1}{(1-z)(1-uz)^{s+1}(1-u)} &= \binom{s+j+1}{s+1}, \quad \text{and} \\ [z^n u^j] \frac{1}{(1-z)^2(1-uz)^2} &= (j+1)(n+1-j). \end{aligned}$$

Theorem 3.2 is an immediate consequence of the formulæ above. □

Corollary 3.1. *The expected value and variance of $D_{n,j}$ are, respectively,*

$$\begin{aligned} \mathbb{E}[D_{n,j}] &= H_j + H_{n+1-j} - 1, \\ \mathbb{V}[D_{n,j}] &= 2(n+1)H_n - (2j+1)H_j - (2n-2j+3)H_{n+1-j} \\ &\quad + 2(n+2) - H_j^2 - H_{n+1-j}^2 - 2H_jH_{n+1-j}. \end{aligned}$$

Furthermore, for $j = \alpha n$, with $0 < \alpha < 1$, we have

$$\begin{aligned} \mathbb{E}[D_{n,\alpha n}] &= 2 \log n + \log \alpha + \log(1-\alpha) + 2\gamma - 1 + o(1), \\ \mathbb{V}[D_{n,\alpha n}] &= 2n \left(1 - \alpha \log \alpha - (1-\alpha) \log(1-\alpha) \right) + \mathcal{O}(\log^2 n), \end{aligned}$$

where $\gamma = 0.5772156649\dots$ is Euler's constant.

To recover higher order ordinary moments, we only need to express the ordinary powers as linear combinations of the shifted falling factorials with coefficients $\lambda_{s,k}$. Thus

$$x^s = \sum_{k=0}^s \lambda_{s,k} (x+2)^{\underline{k}}.$$

It is easy to show that

$$\lambda_{s,k} = \sum_{i=k}^s \left\{ \begin{matrix} i \\ k \end{matrix} \right\} \binom{s}{i} (-2)^{s-i},$$

where $\left\{ \begin{matrix} i \\ k \end{matrix} \right\}$ denote Stirling numbers of the second kind. The coefficients $\lambda_{s,k}$ satisfy a recursion that is similar to that of the Stirling numbers

$$\lambda_{s+1,k} = \lambda_{s,k-1} + (k-2)\lambda_{s,k},$$

and $\lambda_{s,0} = (-2)^s$.

Let us consider now D_n , the number of descendants of a random node in a random BST of size n . The following two theorems give closed formulæ for the probability that D_n is m and for the shifted factorial moments of D_n , i.e. for $d_n^{(s)} = \mathbb{E}[(D_n+2)^{\underline{s}}]$.

Theorem 3.3. *The probability that a randomly chosen internal node in a random binary search tree of size n has m descendants is given by*

$$\mathbb{P}[D_n = m] = \begin{cases} \frac{2(n+1)}{n(m+1)(m+2)} & \text{for } 1 \leq m < n, \\ \frac{1}{n} & \text{for } m = n. \end{cases}$$

Proof. Plug $u = 1$ into the solution of (7) to get

$$\overline{D}(z, v) = D(z, 1, v) = -\frac{2(1-v)}{v^2(1-z)^2} \log \frac{1}{1-vz} - \frac{z(zv - v^2 + 2v - 2)}{v(1-vz)(1-z)^2}. \tag{8}$$

The coefficient of $[z^n v^m] \overline{D}(z, v)$, divided by n , is the sought probability. □

Theorem 3.4. *The s^{th} shifted factorial moment $d_n^{(s)} = \mathbb{E}[(D_n + 2)^s]$ of the number of descendants of a random node in a random binary search tree of size n is given by*

1. $d_n = d_n^{(1)} = 2(1 + \frac{1}{n})H_n - 1$,
2. $d_n^{(2)} = 3(n + 1)$.
3. For all $s \geq 3$,

$$d_n^{(s)} = \frac{1}{n} \left((n+2)^s + \frac{2}{s-1}(n+1)^s \right) \sim \frac{s+1}{s-1} n^{s-1}.$$

Proof. Repeated differentiation of the generating function $v^2 \overline{D}(z, v)$ w.r.t. v and setting $v = 1$, gives us the generating functions of the shifted factorial moments. Their coefficients are extracted much in the same way as in Theorem 3.2. □

A few comments concerning the last theorem are in order now. Observe that for $s \geq 3$

$$\frac{1}{n} \sum_{j=1}^n d_{n,j}^{(s)} = \frac{(n+1)^{s-1}}{s-1} \left(s + 1 + \frac{2}{n} \right).$$

Asymptotically, this quantity is

$$\sim \frac{s+1}{s-1} n^{s-1},$$

one of the observations in the work of Lent [21]. The coincidence in asymptotic behavior with $d_n^{(s)}$ is remarkable; recall that in general

$$\mathbb{E}[D_n^s] \neq \mathbb{E} \left[\left(\frac{1}{n} \sum_{1 \leq j \leq n} D_{n,j} \right)^s \right],$$

except when $s = 1$ and the same observation holds for the shifted factorial moments we were dealing with.

Last, but not least, we can obtain the following corollaries, from Propositions 1.4 and 1.5 and the theorems in this section. These results can already be found in [20], although there is a slight difference in $\mathbb{E}[C_n^{(b)}]$, because $n + 1$ comparisons per partition are counted there, while we count $n - 1$ comparison per partition.

Corollary 3.2. *The expected number of pages in a random binary search tree of size n with page capacity b is*

$$\mathbb{E}[\mathcal{P}_n^{(b)}] = 2 \frac{n+1}{b+2}.$$

The filling ratio for binary search trees is thus

$$\gamma_b = \frac{n/b}{\mathbb{E} \left[\mathcal{P}_n^{(b)} \right]} \sim \frac{1}{2}.$$

Corollary 3.3. *The expected number of recursive calls to sort a random permutation of size n , when the recursion stops in subfiles of size $\leq b$ is*

$$\mathbb{E} \left[R_n^{(b)} \right] = \frac{2n - b}{b + 2}.$$

Also, the expected number of comparisons to sort a random permutation of size n , when the recursion stops in subfiles of size $\leq b$ is

$$\mathbb{E} \left[C_n^{(b)} \right] = 2(n + 1) (H_n - H_{b+1}) + n + 5 - \frac{6(n + 1)}{b + 2}.$$

4. THE NUMBER OF ASCENDANTS IN RANDOM BSTs

Considering the element k of the root of a BST, we obtain for the number of ascendants $A_{n,j}$ of the j^{th} node of a BST of size n the following recursion:

$$\mathbb{P} [A_{n,j} = m] = \frac{1}{n} \llbracket m = 1 \rrbracket + \frac{1}{n} \sum_{k=1}^{j-1} \mathbb{P} [A_{n-k,j-k} = m - 1] + \frac{1}{n} \sum_{k=j+1}^n \mathbb{P} [A_{k-1,j} = m - 1]. \quad (9)$$

Introducing the generating function for the family of random variables $\{A_{n,j}\}$

$$A(z, u, v) = \sum_{1 \leq j, m \leq n} \mathbb{P} [A_{n,j} = m] z^n u^j v^m,$$

this recursion translates by multiplying both sides by $nz^{n-1}u^jv^m$ and summing for all $n \geq 1$, $1 \leq j \leq n$ and $m \geq 1$ into the following differential equation:

$$\frac{\partial A}{\partial z} = \frac{v}{1 - z} A + \frac{uv}{1 - uz} A + \frac{uv}{(1 - z)(1 - uz)}$$

with the initial condition $A(0, u, v) = 0$. This differential equation has the following solution

$$A(z, u, v) = \frac{uv}{(1 - z)^v(1 - uz)^v} \int_0^z (1 - t)^{v-1}(1 - ut)^{v-1} dt. \quad (10)$$

Starting with this generating function, it is easy to get the following theorems. At first we obtain an old result from [3]:

Theorem 4.1. *The expected number of ascendants $a_{n,j} = \mathbb{E} [A_{n,j}]$ of the j^{th} node in a random binary search tree of size n is*

$$a_{n,j} = H_j + H_{n+1-j} - 1.$$

Proof. Starting with (10), taking derivatives w.r.t. v and setting $v = 1$, we get the generating function $\mathcal{A}(z, u)$, whose coefficients are the expected values $a_{n,j} = \mathbb{E} [A_{n,j}]$. It is given by

$$\mathcal{A}(z, u) = \frac{u}{(1 - z)(1 - uz)} \log \frac{1}{1 - z} + \frac{1}{(1 - z)(1 - uz)} \log \frac{1}{1 - uz} - \frac{uz}{(1 - z)(1 - uz)}.$$

It is easy to extract the coefficients of this expression, which leads immediately to the stated theorem. \square

Theorem 4.2. *The second factorial moment $a_{n,j}^{(2)} = \mathbb{E} \left[(A_{n,j})^2 \right]$ of the number of ascendants of the j^{th} node in a random binary search tree of size n is*

$$a_{n,j}^{(2)} = \frac{2(n+1)}{(n+1-j)j} H_n + H_j^2 + 2H_j H_{n+1-j} + \frac{2(-nj - n + j^2 - j - 1)}{(n+1-j)j} H_j + H_{n+1-j}^2 + \frac{2(-nj - n + j^2 - j - 1)}{(n+1-j)j} H_{n+1-j} - H_j^{(2)} - H_{n+1-j}^{(2)} - \frac{2(-2nj + 2j^2 - 2j - 1)}{(n+1-j)j}. \tag{11}$$

Proof. Differentiating equation (10) two times w.r.t. v and setting $v = 1$ gives the generating function $\mathcal{A}^{(2)}(z, u)$ of the second factorial moments $a_{n,j}^{(2)}$ of the number of ascendants:

$$\begin{aligned} \mathcal{A}^{(2)}(z, u) &= -\frac{2zu}{(1-uz)(1-z)} \log \frac{1}{1-uz} - \frac{2zu}{(1-uz)(1-z)} \log \frac{1}{1-z} \\ &\quad - \frac{2(uz - u - 1)}{(1-uz)(1-z)} \log \frac{1}{1-z} \log \frac{1}{1-uz} + \frac{u}{(1-uz)(1-z)} \log^2 \frac{1}{1-z} \\ &\quad + \frac{1}{(1-uz)(1-z)} \log^2 \frac{1}{1-uz} + \frac{2u}{(1-uz)(1-z)} \int_0^z \log \frac{1}{1-t} \log \frac{1}{1-ut} dt. \end{aligned}$$

Extracting the coefficients leads to the given theorem. Since one expression in $\mathcal{A}^{(2)}(z, u)$ turns out to be a bit messier, we sketch how to extract the coefficients of it. First we get the following sum

$$\begin{aligned} [z^n u^j] \frac{1}{(1-z)(1-uz)} \int_0^z \log \frac{1}{1-t} \log \frac{1}{1-ut} dt &= \sum_{k=0}^j \sum_{l=0}^{n-j+k} [z^l u^k] \int_0^z \log \frac{1}{1-t} \log \frac{1}{1-ut} dt \\ &= \sum_{k=1}^j \sum_{l=k+2}^{n-j+k} \frac{1}{lk(l-k-1)}, \end{aligned}$$

which can be simplified to

$$\begin{aligned} \sum_{k=1}^j \sum_{l=k+2}^{n-j+k} \frac{1}{lk(l-k-1)} &= \sum_{k=1}^j \frac{1}{k} \sum_{l=1}^{n-j-1} \frac{1}{l(l+k+1)} = \sum_{k=1}^j \frac{1}{k(k+1)} \sum_{l=1}^{n-j-1} \left(\frac{1}{l} - \frac{1}{l+k+1} \right) \\ &= \sum_{k=1}^j \frac{1}{k(k+1)} (H_{n-j-1} + H_{k+1} - H_{n-j+k}) = \sum_{k=1}^j \left(\frac{1}{k} - \frac{1}{k+1} \right) (H_{n-j-1} + H_{k+1} - H_{n-j+k}) \\ &= H_{n-j-1} \sum_{k=1}^j \left(\frac{1}{k} - \frac{1}{k+1} \right) + \sum_{k=1}^j \left(\frac{H_k}{k} - \frac{H_{k+1}}{k+1} \right) + \sum_{k=1}^j \left(\frac{1}{k} - \frac{1}{k+1} \right) \\ &\quad - \sum_{k=1}^j \left(\frac{H_{n-j+k}}{k} - \frac{H_{n-j+k+1}}{k+1} \right) - \frac{1}{n-j} \sum_{k=1}^j \left(\frac{1}{k+1} - \frac{1}{n-j+k+1} \right). \end{aligned}$$

The sums telescope and we finally get

$$\begin{aligned} [z^n u^j] \frac{1}{(1-z)(1-uz)} \int_0^z \log \frac{1}{1-t} \log \frac{1}{1-ut} dt &= \frac{n+1}{(j+1)(n-j)} (H_{n+1} - H_{j+1} - H_{n+1-j}) \\ &\quad + \frac{2jn^2 - 4nj^2 + 2j^3 + n^2 - jn + 2n - 2j + 1}{(n-j)(j+1)(n+1-j)}. \end{aligned}$$

□

The next theorem gives the variance, which is now easy to obtain.

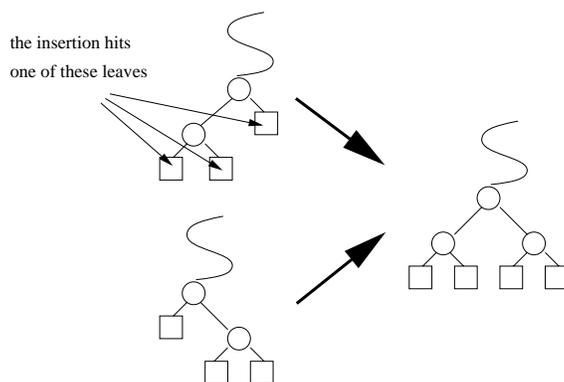


FIGURE 2. The fringe heuristic

Theorem 4.3. *The variance $\mathbb{V}[A_{n,j}]$ of the number of ascendants of the j^{th} node in a random search tree of size n is*

$$\begin{aligned} \mathbb{V}[A_{n,j}] = & \frac{2(n+1)}{(n+1-j)j} H_n + \frac{(nj - 2n - j^2 + j - 2)}{(n+1-j)j} H_j + \frac{(nj - 2n - j^2 + j - 2)}{(n+1-j)j} H_{n+1-j} \\ & - H_j^{(2)} - H_{n+1-j}^{(2)} + \frac{2(nj - j^2 + j + 1)}{(n+1-j)j}. \end{aligned} \quad (12)$$

□

5. LOCALLY BALANCED BINARY SEARCH TREES

One approach to avoid drastically unbalanced binary search trees is the introduction of strict balance constraints like in AVLs or red-black trees [1, 11]. Such schemes guarantee logarithmic performance of searches and updates in the worst-case, but they have additional space requirements and are more difficult to implement than standard BSTs. As an alternative, several authors [4, 34, 27] have suggested the use of a simple heuristic that makes the construction of poorly balanced trees much less likely than with the use of the standard algorithms. Furthermore, the heuristic was shown to yield significant savings in the expected search time.

The basic idea is really simple: whenever a son is appended to a node that itself is a single son (its “brother” is an external node), a rotation of the three nodes is performed to place the median of the three elements as the root of the subtree and the other two elements as sons (see Figure 2). Since no other kind of rebalancing operation is ever made, Poblete and Munro refer to this technique as a *fringe* heuristic. We will call the binary search trees constructed in this way *locally balanced binary search trees* (LBST, for short).

Poblete and Munro [27] and Poblete [26] carry on the analysis of this heuristic and some generalizations by means of bottom-up or fringe techniques: they basically study the number of nodes that are at level k and which are the root of a subtree of size 1 or 2.

As we have already mentioned in the introduction, the standard model for random LBSTs states that a random LBST of size n is the result of n random insertions into an initially empty tree. Equivalently, a random LBST of size n is the result of inserting the elements of a random permutation of $\{1, \dots, n\}$ into an initially empty tree. Here, we show that a recursive, top-down definition of the randomness model is also possible. This characterization of the model of randomness is more amenable to the kind of algebraic manipulations that we want to carry on; as we will see, the recurrence relations for the analyzed quantities translate to equations over generating functions in a natural way, almost automatically.

Definition 5.1. 1. A random binary search tree of size $s \leq 2$ is also a random locally balanced search tree. Recall that a BST of size 2 is random if the smallest (resp. largest) key is the root with probability $1/2$.

2. A binary search tree T of size $n \geq 3$, with left and right subtrees T_1 and T_2 , is a random LBST if and only if, both T_1 and T_2 are random independent LBSTs, and

$$\pi_{n,k} = \mathbb{P} \left[|T_1| = k - 1 \mid |T| = n \right] = \frac{(k - 1)(n - k)}{\binom{n}{3}}, \quad \text{for all } 1 \leq k \leq n.$$

The reader should have noticed that the only difference between this definition and that for random BSTs relies on the *splitting probabilities* $\pi_{n,k}$. In the case of BSTs, each element of the random permutation has the same probability (namely, $1/n$) of being the first element and hence of becoming the root. In the case of LBSTs, when $n \geq 3$, the probability that the k^{th} element is one of the first three elements of the permutation and is the median of these three elements is

$$\frac{1}{n} \times \frac{k - 1}{n - 1} \times \frac{n - k}{n - 2} \times 3! = \pi_{n,k}.$$

Indeed, the left hand side of the equation above give us the probability that the k^{th} element is the first, times the probability that it is followed by a smaller element, times the probability that the two elements are followed by a larger element. For any permutation of such three elements, we have that the k^{th} element is among the first three elements and it is their median. Now, under these conditions the k^{th} element will be the root of the LBST (after the insertion of the first three elements, with rebalancing if necessary). The insertion of the fourth, fifth, etc. elements will not affect the root of the LBST. The principle applies recursively to the subsequences of elements smaller and greater than the selected element and the definition follows.

This argument also justifies the deep connection between LBSTs, quicksort and quickselect (see Propositions 1.3 and 1.5), when we consider the variants that select the median of 3 elements taken at random as the pivot of each partitioning phase.

6. THE NUMBER OF DESCENDANTS IN RANDOM LBSTs

As in Section 3, let $D_{n,j}$ denote the number of descendants of the j^{th} node, but now in a random LBST of size n . The recursion for $\mathbb{P} [D_{n,j} = m]$ is almost the same as for random BSTs, the only difference being the splitting probability $\pi_{n,k}$, the probability that the root of the LBST is the k^{th} element. Thus,

$$\begin{aligned} \mathbb{P} [D_{n,j} = m] &= \left[\sum_{1 \leq k < j} \pi_{n,k} \mathbb{P} [D_{n-k,j-k} = m] + \pi_{n,j} \mathbb{P} [m = n] \right. \\ &\quad \left. + \sum_{j < k \leq n} \pi_{n,k} \mathbb{P} [D_{k-1,j} = m] \right]. \end{aligned} \tag{13}$$

Theorem 6.1. *Let*

$$D_z(z, u, v) = \frac{\partial}{\partial z} \sum_{n,j,m} \mathbb{P} [D_{n,j} = m] z^n u^j v^m = \sum_{n,j,m} n \mathbb{P} [D_{n,j} = m] z^{n-1} u^j v^m.$$

Then,

$$\begin{aligned} D_z(z, u, v) &= \frac{A_0(z, u, v)}{v(1 - z)^2(1 - uz)^2(1 - uv)^2(1 - v)^2(v - u)^2(1 - u)^2} \\ &\quad + \frac{A_1(z, u, v)}{(1 - uz)^2(1 - u)(1 - v)^3(1 - uv)^3} \log \frac{1}{1 - z} \\ &\quad + \frac{A_2(z, u, v)}{(1 - z)^2(u - v)^3(1 - u)(1 - v)^3} \log \frac{1}{1 - uz} \end{aligned}$$

$$\begin{aligned}
 &+ \frac{A_3(z, u, v)}{v^2(1-z)^2(1-uz)^2(v-u)^3(1-u)^3(1-v)^3} \log \frac{1}{1-vz} \\
 &+ \frac{A_4(z, u, v)}{v^2(1-z)^2(1-uz)^2(1-u)^3(1-v)^3(1-uv)^3} \log \frac{1}{1-uvz},
 \end{aligned}$$

where each of the $A_i(z, u, v)$'s is a complicated polynomial in z, u and v . They are listed in full in the appendix.

Proof. We multiply the recursion (13) by $\binom{n}{3}$ and $z^{n-3}u^jv^m$, sum up over all $n \geq 1$ and $1 \leq j, m \leq n$ to get the following differential equation:

$$\frac{1}{6} \frac{\partial^2 D_z}{\partial z^2} = \frac{u^2 v^3}{(1-vz)^2(1-uvz)^2} + \frac{u^2}{(1-uz)^2} D_z + \frac{1}{(1-z)^2} D_z, \tag{14}$$

where the initial conditions are $D_z(0, u, v) = uv$ and $\frac{\partial}{\partial z} D_z(0, u, v) = uv(1+u)(1+v)$. We use the partial derivative w.r.t. z to define $D_z(z, u, v)$ because the differential equation just given, which translates the recurrence for $\mathbb{P}[D_{n,j} = m]$, is then of the second order. Had we introduced the generating function $D_z(z, u, v)$ in the standard manner, we would have had a third order differential equation, with no appearance of the function itself, only the first and third derivatives.

The differential equation (14) is solvable: its explicit form (abridged) is the one given in the statement of the theorem. \square

From the explicit form of $D_z(z, u, v)$ given in Theorem 6.1 we can, in principle, compute exact expressions for $\mathbb{P}[D_{n,j} = m]$ and all moments. However, the task is daunting, and we will content ourselves computing the expected value and the second factorial moment in the next two theorems.

Theorem 6.2. *The expected number of descendants $d_{n,j} = \mathbb{E}[D_{n,j}]$ of the j^{th} node in a random LBST of size n is, when $5 \leq j \leq n - 4$*

$$\begin{aligned}
 d_{n,j} = &-\frac{12}{7}H_n + \frac{12}{7}H_j + \frac{12}{7}H_{n+1-j} \\
 &-\frac{6}{7j} - \frac{6}{7(n+1-j)} + \frac{79}{70} \\
 &-\frac{3(3j-5)}{7n} + \frac{6(j-1)^2}{7n^2} + \frac{2(2j-3)(j-1)^2}{7n^3} \\
 &+ \frac{3(j-2)(j-1)^3}{7n^4} - \frac{3(2j-5)(j-1)^4}{7n^5} + \frac{2(j-3)(j-1)^5}{7n^6}.
 \end{aligned}$$

The remaining cases when $j \leq 4$ (or when $j > n - 4$, by symmetry) appear in the appendix.

Proof. Taking the first derivative with respect to v , and setting $v = 1$ we get⁴

$$\begin{aligned}
 \left. \frac{\partial D_z}{\partial v} \right|_{v=1} &= \frac{B_0(z, u)}{70(1-uz)^2(1-u)^7(1-z)^2} \\
 &+ \frac{B_1(z, u)}{7(1-u)^7(1-uz)^2} \log \frac{1}{1-z} \\
 &+ \frac{B_2(z, u)}{7(1-z)^2(1-u)^7} \log \frac{1}{1-uz},
 \end{aligned}$$

where the $B_i(z, u)$'s are polynomials in z and u . Their explicit value can be found in the appendix at the end of this paper.

⁴It turns out that Maple gets stuck doing the work in the obvious way, i.e. take the derivative, then take the limit when $v \rightarrow 1$. But we can produce the differential equations satisfied by the generating functions for the factorial moments from the differential equation (14) and solve them. Also, the problem can be fixed by computing a series expansion of the derivatives around $v = 1$.

In order to get to the coefficients, we use formulæ such as

$$\begin{aligned}
 [z^n u^j] \frac{1}{(1-u)^4(1-uz)^2} \log \frac{1}{1-z} &= (n+1) \binom{j-n+3}{3} (H_n - H_{n-1-j}) \\
 &+ \frac{(j+1)n^3}{6} - \frac{5(j+1)(j+2)n^2}{12} \\
 &+ \frac{(j+1)(11j^2 + 40j + 30)n}{36} - \frac{3j-10}{12} \binom{j+3}{3},
 \end{aligned}$$

$$\begin{aligned}
 [z^n u^j] \frac{1}{(1-u)^5(1-uz)^2} \log \frac{1}{1-z} &= (n+1) \binom{j-n+4}{4} (H_n - H_{n-1-j}) \\
 &- \frac{(j+1)n^4}{24} + \frac{(7j+18)(j+1)n^3}{48} \\
 &- \frac{(j+1)(13j^2 + 65j + 75)n^2}{72} \\
 &+ \frac{(j+1)(25j^3 + 173j^2 + 348j + 180)n}{288} \\
 &- \frac{12j-65}{60} \binom{j+4}{4},
 \end{aligned}$$

and similar ones that are not too hard to obtain. To retrieve the final answer, we have also to take into account that we need to shift the coefficients in z^n by 1 and multiply by $\frac{1}{n}$, because we were considering $\frac{\partial}{\partial z} \sum_{n,j,m} \mathbb{P}[D_{n,j} = m] z^n u^j v^m$. Putting everything together, the theorem follows. \square

Theorem 6.3. *The second factorial moment of the number of descendants $d_{n,j}^{(2)} = \mathbb{E} [D_{n,j}^2]$ of the j^{th} node in a random LBST of size n is, when $5 \leq j \leq n-4$,*

$$\begin{aligned}
 d_{n,j}^{(2)} &= \left(\frac{36n}{5} - \frac{12}{35} \right) H_n + \left(\frac{36j}{5} - \frac{36n}{5} - \frac{48}{7} \right) H_{n+1-j} + \left(\frac{12}{35} - \frac{36j}{5} \right) H_j \\
 &- \frac{132}{35j} - \frac{132}{35(n+1-j)} + \frac{3489}{175} - \frac{33j}{5} + \left(\frac{66}{7} - \frac{429j}{35} + \frac{33j^2}{5} \right) \frac{1}{n} \\
 &+ \frac{132(j-1)^2}{35n^2} + \frac{44(2j-3)(j-1)^2}{35n^3} + \frac{66(j-2)(j-1)^3}{35n^4} \\
 &- \frac{66(2j-5)(j-1)^4}{35n^5} + \frac{44(j-3)(j-1)^5}{35n^6}.
 \end{aligned}$$

The formulæ for the second factorial moment in the special cases (when $j \leq 4$ or $j > n-4$) are collected in a table in the appendix.

Proof. The second factorial moment $d_{n,j}^{(2)}$ is the coefficient of $z^{n-1}u^j$ times $\frac{1}{n}$ in

$$\begin{aligned}
 \frac{\partial^2 D_z}{\partial v^2} \Big|_{v=1} &= \frac{C_0(z, u)}{35(1-z)^2(1-u)^7(1-uz)^2} \\
 &+ \frac{C_1(z, u)}{35(1-z)^2(1-u)^7(1-uz)^2} \log \frac{1}{1-z} \\
 &+ \frac{C_2(z, u)}{35(1-z)^2(1-u)^7(1-uz)^2} \log \frac{1}{1-uz}
 \end{aligned}$$

where the $C_i(z, u)$'s are polynomials in z and u . They have been listed in the appendix. Using techniques similar to the ones in the proof of Theorem 6.2, we extract the coefficients and obtain the stated result. \square

As in Section 3, we shift now our attention to the number of descendants of a random node in a random LBST of size n . We start giving an explicit expression for the probability distribution of D_n .

Theorem 6.4. *The probability that a random node in a random LBST of size n has m descendants is*

$$\mathbb{P}[D_n = m] = \frac{12}{7} \frac{(n+2+m)(n-1-m)}{n^2(m+1)(m+2)} - \frac{12}{7} \frac{m^{\underline{5}}}{nn^{\underline{6}}} + \frac{12}{7n^2},$$

for $5 \leq m < n$. The probability that a random node has n descendants is $\mathbb{P}[D_n = n] = \frac{1}{n}$. Furthermore, the probability that a random node in a random LBST of size n has no children is

$$\mathbb{P}[D_n = 1] = \frac{6}{7} \frac{1}{n^2} \binom{n+1}{2} = \frac{3}{7} \left(1 + \frac{1}{n}\right).$$

In the appendix, a table collects the general result for $5 \leq m < n$ as well as the special cases where $m < 5$ or $m = n$.

Proof. If we consider the explicit form for $D_z(z, u, v)$ given in Theorem 6.1 and average w.r.t. j , i.e. we plug $u = 1$ there, we get

$$\begin{aligned} \bar{D}_z(z, v) &= \frac{v}{(1-vz)^2} + \frac{12}{7(1-v)(1-z)} - \frac{24}{7v(1-z)^2} + \frac{2(v^2 - 6v + 12)}{7v(1-z)^3} \\ &\quad + \frac{v}{7(1-v)^5} \left[-15(1-v)^3 + 20v(1-v)^2(1-z) - 30v^2(1-v)(1-z)^2 \right. \\ &\quad \left. + 60v^3(1-z)^3 + (1-7v+23v^2-57v^3-22v^4+2v^5)(1-z)^4 \right] \\ &\quad - \frac{60}{7} \frac{v^5(1-z)^4}{(1-v)^6} \log \frac{1}{1-z} + \left(\frac{60}{7} \frac{v^5(1-z)^4}{(1-v)^6} - \frac{24}{7} \frac{1-v}{v^2(1-z)^3} \right) \log \frac{1}{1-vz}. \end{aligned} \quad (15)$$

Alternatively, we can write down the differential equation for $\bar{D}_z(z, v) = D_z(z, 1, v)$ and solve it. The differential equation is

$$\frac{1}{6} \frac{\partial^2 \bar{D}_z}{\partial z^2} = \frac{v^3}{(1-vz)^4} + 2 \frac{\bar{D}_z}{(1-z)^2},$$

where the initial conditions are $\bar{D}_z(0, v) = v$, and $\frac{\partial}{\partial z} \bar{D}_z(0, v) = 2v(1+v)$. The reader may readily check that the explicit form given in Equation (15) is a solution to the differential equation above.

The purely rational term in Equation (15), i.e. the one that is not multiplied by any logarithmic function, although more complicated than the others, has the very pleasant feature that “almost” all coefficients are $\frac{12}{7}$. On the other hand,

$$[z^n v^m] \frac{1}{(1-z)^3} \log \frac{1}{1-vz} = \frac{1}{m} \binom{n-m+2}{2},$$

and thus

$$-[z^n v^m] \frac{24}{7} \frac{1-v}{v^2} \frac{1}{(1-z)^3} \log \frac{1}{1-vz} = \frac{12}{7} \frac{(n+3+m)(n-m)}{(m+2)(m+1)}.$$

This is the main contribution in the coefficient $z^n v^m$ of $\bar{D}_z(z, v)$, the remaining contributions being small. Indeed,

$$\frac{60}{7} \frac{v^5(1-z)^4}{(1-v)^6} \log \frac{1}{1-vz}$$

produces no coefficients at all, since $m \leq n$. And the remaining contribution comes from

$$-[z^n v^m] \frac{60}{7} \frac{v^5(1-z)^4}{(1-v)^6} \log \frac{1}{1-z} = -\frac{12}{7} \frac{m^{\underline{5}}}{n^{\underline{6}}}.$$

The general part of the theorem follows from the considerations made above. The special cases, when $m < 5$ or $m = n$ have to be dealt with separately. In particular, to get the probability that a random node in a random LBST of size n has no children (the special case $m = 1$) we compute

$$\left. \frac{\partial \overline{D}_z(z, v)}{\partial v} \right|_{v=0} = \frac{6}{7} \frac{1}{(1-z)^3} + \frac{1}{7}(1-z)^4,$$

extract the coefficient of z^{n-1} in the GF above and divide by n^2 , yielding $\mathbb{P}[D_n = 1] \sim 3/7$. Also, for evident reasons, $\mathbb{P}[D_n = n] = 1/n$, since only the root has n descendants and we choose it with probability $1/n$. \square

Finally, the moments of D_n can be computed after differentiation of $\overline{D}_z(z, v)$, whose explicit form was given in the proof above. We state now the following result.

Theorem 6.5. *Let $d_n^{(s)} = \mathbb{E}[(D_n + 2)^{\underline{s}}]$, i.e., $d_n^{(s)}$ is the shifted s^{th} factorial moment of the number of descendants of a random node in a random locally balanced binary search tree of size n . Furthermore, let $d_n = d_n^{(1)}$. Then*

1. $d_n = \frac{12}{7} \left(1 + \frac{1}{n}\right) H_n - \frac{1}{49} \left(26 - \frac{9}{n}\right)$, for $n \geq 6$,
2. $d_n^{(2)} = \frac{5(n+1)(7n+2)}{14n}$, for $n \geq 6$,
3. $d_n^{(3)} = \frac{(n+1)(10n^2+5n+6)}{6n}$, for $n \geq 6$.
4. For all $n \geq s + 7$ and all $s \geq 4$,

$$d_n^{(s)} = \frac{A(s, n) (n+1)^{\underline{s+1}}}{(s+6)^{\underline{6}} (n+2-s)^{\underline{2}} n n^{\underline{6}} (s-1)},$$

where

$$\begin{aligned} A(s, n) = & (s+5)^{\underline{5}}(s+3)(s+2)n^7 \\ & - (s+4)^{\underline{4}}(s+2)(13s^2+128s+195)n^6 \\ & + (s+3)^{\underline{3}}(67s^4+1082s^3+6125s^2+11326s+6600)n^5 \\ & - 5(s+2)^{\underline{2}}(35s^5+643s^4+4459s^3+15317s^2+15906s+3960)n^4 \\ & + 4(s+1)(61s^6+1159s^5+8157s^4+24383s^3+60116s^2-9276s-31680)n^3 \\ & - 4(43s^7+794s^6+5176s^5+10190s^4-80183s^3+29336s^2-220956s-77040)n^2 \\ & + 48(s^7+17s^6+97s^5+215s^4+1894s^3-39832s^2+41208s-25200)n \\ & - 1036800s^2+3110400s-2073600. \end{aligned}$$

Corollary 6.1. *For any $n \geq 6$ and for $j = \alpha n$, with $0 < \alpha < 1$, we have*

$$\begin{aligned} \mathbb{E}[D_{n,\alpha n}] &= \frac{12}{7} \log n + \mathcal{O}(1), \\ \mathbb{V}[D_{n,\alpha n}] &= -\frac{3}{5}n \left(11\alpha(1-\alpha) + 12\alpha \log \alpha + 12(1-\alpha) \log(1-\alpha)\right) + \mathcal{O}(\log^2 n). \end{aligned}$$

As in Section 3, several interesting corollaries may be deduced from the results in this section and Propositions 1.4 and 1.5.

Corollary 6.2. *The expected number of pages in a random locally balanced search tree of size n with page capacity $b \geq 2$ is*

$$\mathbb{E} \left[\mathcal{P}_n^{(b)} \right] \sim \frac{12}{7} \frac{n}{b+2}.$$

The filling ratio for locally balanced search trees is thus

$$\gamma_b = \frac{n/b}{\mathbb{E} \left[\mathcal{P}_n^{(b)} \right]} \sim \frac{7}{12} = 0.58333 \dots$$

Corollary 6.3. *The expected number of recursive calls to sort a random permutation of size n , when the recursion stops at subfiles of size $\leq b$ and the pivots are selected as the median of samples of three elements, is*

$$\mathbb{E} \left[R_n^{(b)} \right] = \mathbb{E} \left[\mathcal{P}_n^{(b)} \right] - 1 \sim \frac{12}{7} \frac{n}{b+2}.$$

Also of interest is the expectation $C_{n,b} := \mathbb{E} \left[C_n^{(b)} \right]$ of the number of comparisons to sort a random permutation of size n with quicksort, where the pivots are selected as the median of samples of three elements (for subfiles of length $n \geq 3$) and the recursion stops at subfiles of size $\leq b$. We only consider here comparisons, that appear by comparing the pivot to each other element in the partitioning step, and do not count the (on average) $\frac{8}{3}$ comparisons to select the median of three elements. We also make the assumption, that small subfiles of size $n \leq b$ are stored unsorted in own pages and so we do not count comparisons in these cases. To get these expectations we don't use Proposition 1.5. We take another approach and start with the following recursion for $C_{n,b}$:

$$C_{n,b} = n - 1 + \sum_{k=1}^n \pi_{n,k} (C_{k-1,b} + C_{n-k,b}) \quad \text{for } n > b \geq 0 \text{ and } n \geq 3, \tag{16}$$

with initial values $C_{2,0} = 1$, $C_{2,1} = 1$ and $C_{n,b} = 0$ otherwise. (With these initial values we take care of the one additional comparison, sorting a subfile of length 2, when the pages are smaller than 2.)

To solve this recurrence, we introduce the bivariate generating function $C_z(z, v) = \sum_{n>b \geq 0} C_{n,b} n z^{n-1} v^b$. Multiplying both sides of equation (16) by $n(n-1)(n-2)z^{n-3}v^b$ and summing up over all $n > b \geq 0$ leads to the following differential equation

$$\frac{\partial^2}{\partial z^2} C_z(z, v) = \frac{12}{(1-z)^2} C_z(z, v) + \frac{12(z^6 v^4 + z^5 v^4 + z^5 v^3 - 15 z^4 v^3 + 10 z^3 v^3 + 10 z^3 v^2 - 5 z^2 v^3 + 5 z^2 v^2 - 5 z^2 v + z v^3 - 4 z v^2 - 4 z v + z + v^2 + v + 1)}{(1-z)^5 (1-zv)^5}, \tag{17}$$

with initial conditions $C_z(0, v) = 0$ and $\frac{\partial}{\partial z} C_z(0, v) = 2(1+v)$.

This differential equation is of Eulerian type, and can be solved easily. We get then

$$\begin{aligned}
 C_z(z, v) = & \left(\frac{120}{7} \frac{(1-z)^4(v+2)v^5}{(1-v)^8} + \frac{24}{7} \frac{1}{(1-v)(1-z)^3} \right) \log \frac{1}{1-z} \\
 & + \left(-\frac{120}{7} \frac{(1-z)^4(v+2)v^5}{(1-v)^8} + \frac{24}{7} \frac{2v-3}{(1-v)(1-z)^3v^2} \right) \log \frac{1}{1-zv} \\
 & - 12 \frac{v}{(1-v)^3(1-zv)} + 2 \frac{v}{(1-zv)^2(1-v)} - 2 \frac{v}{(1-zv)^3(1-v)} \\
 & - \frac{2}{49} \frac{89v-252}{(1-v)(1-z)^3v} - \frac{2}{7} \frac{7v^2-31v+36}{(1-z)^2(1-v)^2v} - \frac{12}{7} \frac{2v-3}{(1-v)^3(1-z)} \\
 & + \frac{2}{49} \frac{R(z, v)}{(1-v)^7}.
 \end{aligned} \tag{18}$$

with

$$\begin{aligned}
 R(z, v) = & 40z^4v^6 + 929z^4v^5 + 327z^4v^4 - 23z^4v^3 - 23z^4v^2 + 12z^4v - 2z^4 - 160z^3v^6 - 3296z^3v^5 - \\
 & 468z^3v^4 + 92z^3v^3 + 92z^3v^2 - 48z^3v + 8z^3 + 240z^2v^6 + 4104z^2v^5 - 768z^2v^4 + 282z^2v^3 - 138z^2v^2 + \\
 & 72z^2v - 12z^2 - 160zv^6 - 1896zv^5 + 1632zv^4 - 1168zv^3 + 372zv^2 - 48zv + 8z + 40v^6 + 54v^5 - \\
 & 618v^4 + 1132v^3 - 828v^2 + 222v - 2.
 \end{aligned}$$

Extracting the coefficients, we get with $\mathbb{E} [C_n^{(b)}] = C_{n,b} = \frac{1}{n} [z^{n-1}v^b]C_z(z, v)$ the required expectations. This leads to

Theorem 6.6. *The expected number of comparisons to sort a random permutation of size n , when the recursion stops in subfiles of size $\leq b$ and the pivots are selected as the median of samples of three elements, is for $n > b \geq 0$ and $n \geq 6$ given as*

$$\mathbb{E} [C_n^{(b)}] = \frac{12}{7}(n+1)H_n - \frac{12}{7}(n+1)H_{b+1} + \frac{37n}{49} + \frac{219}{49} - \frac{36(n+1)}{7(b+2)} + \frac{4(3b-1)(b+1)^6}{49n^6}.$$

7. THE NUMBER OF ASCENDANTS OF A GIVEN NODE IN A LBST

As in the case of the number of ascendants in a random BST, computing the probability that the j^{th} node in a random LBST has m ascendants turns out to be an extremely difficult problem.

However, the recursive definition can easily be translated to a differential equation for the corresponding generating function $A_z(z, u, v)$. Because of the same technical reason discussed in Section 6, the function $A_z(z, u, v)$ is actually the derivative w.r.t. z of the generating function such that the coefficient of $z^n u^j$ is the PGF of $A_{n,j}$. The recurrence for $A_{n,j}$

$$A_{n,j} = \sum_{k=1}^{j-1} \pi_{n,k}(A_{n-k,j-k} + 1) + \pi_{n,j} + \sum_{k=j+1}^n \pi_{n,k}(A_{k-1,j} + 1) \quad \text{for } n \geq 3,$$

with initial values $A_{0,j} = 0$, $A_{1,1} = 1$, $A_{2,1} = \frac{3}{2}$, $A_{2,2} = \frac{3}{2}$ and $A_{n,j} = 0$ otherwise, translates into the second-order differential equation

$$\frac{1}{6} \frac{\partial^2 A_z}{\partial z^2} = \frac{v}{(1-z)^2} A_z + \frac{u^2 v}{(1-uz)^2} A_z + \frac{u^2 v}{(1-z)^2(1-uz)^2}, \tag{19}$$

and the initial values are $A_z(0, u, v) = uv$ and $\frac{\partial}{\partial z} A_z(0, u, v) = uv(1+v)(1+u)$. This differential equation is the starting point for our next theorems.

Theorem 7.1. *The expected number of ascendants $a_{n,j} = \mathbb{E}[A_{n,j}]$ of the j^{th} node in a random locally balanced search tree of size n is*

$$\begin{aligned}
 a_{n,j} = & \frac{24}{35}H_n + \frac{18}{35}H_j + \frac{18}{35}H_{n+1-j} \\
 & + \frac{12}{35j} + \frac{12}{35(n+1-j)} - \frac{279}{175} - \frac{6}{7n} \\
 & + \frac{18j}{35n} - \frac{12(j-1)^2}{35n^2} - \frac{4(2j-3)(j-1)^2}{35n^3} \\
 & - \frac{6(j-2)(j-1)^3}{35n^4} + \frac{6(2j-5)(j-1)^4}{35n^5} - \frac{4(j-3)(j-1)^5}{35n^6},
 \end{aligned}$$

for $5 \leq j \leq n-4$. In the appendix we give also the cases $j = 1, 2, 3, 4$. The cases where $j > n-4$ follow from the special cases with $j \leq 4$ and the symmetry in j and $n+1-j$ of $a_{n,j}$.

Proof. Although it is in principle possible to solve the differential equation⁵ (19), it is sufficient for our purpose to take derivatives w.r.t. v and setting $v = 1$, to get the differential equation for $\mathcal{A}_z(z, u)$, the generating function whose coefficients are the expected values $a_{n,j} = \mathbb{E}[A_{n,j}]$. It is

$$\frac{1}{6} \frac{\partial^2 \mathcal{A}_z}{\partial z^2} - \left(\frac{1}{(1-z)^2} + \frac{u^2}{(1-uz)^2} \right) \mathcal{A}_z = \frac{u}{1-u} \left(\frac{1}{(1-z)^4} - \frac{u^3}{(1-uz)^4} \right), \tag{20}$$

and the initial conditions are now $\mathcal{A}_z(0, u) = u$ and $\frac{\partial}{\partial z} \mathcal{A}_z(0, u) = 3u(1+u)$.

The solution of the differential equation (20) yields the explicit form

$$\begin{aligned}
 \mathcal{A}_z(z, u) = & \frac{D_0(z, u)}{(1-z)^2(1-u)^7(1-uz)^2} \\
 & + \frac{D_1(z, u)}{(1-z)^2(1-u)^7(1-uz)^2} \log \frac{1}{1-z} \\
 & + \frac{D_2(z, u)}{(1-z)^2(1-u)^7(1-uz)^2} \log \frac{1}{1-uz},
 \end{aligned}$$

where the polynomials $D_i(z, u)$ can be found in the appendix.

Once we have the explicit form for \mathcal{A}_z , extracting the coefficients is just a matter of patience and careful computations. A possible shortcut is to expand each of the three main parts of \mathcal{A}_z as power series in z and u , and spot a pattern in the shape of the coefficients. The inspired guesses can be readily checked and proved by induction. For instance, the coefficient of $z^n u^j$ in the purely rational term of \mathcal{A}_z is

$$-\frac{69}{175}n + \frac{18}{35}j - \frac{9}{175},$$

whenever $5 \leq j \leq n-4$; the remaining values of j are special cases that we have to consider separately. Similarly, the coefficient of $z^n u^j$ in the second term—the one that contains $\log(1/(1-uz))$ as a factor—is

$$\frac{18H_j(n+1)}{35} - \frac{18j}{35} + \frac{12n}{35j} - \frac{12}{35} + \frac{12}{35j}.$$

In the same vein, an explicit formula for the coefficient of the first term can be obtained. Finally, we collect everything, consider the coefficient $z^{n-1} u^j$ and divide by n , since \mathcal{A}_z is a derivative w.r.t. z . □

⁵With the substitutions $A_z(z, u, v) = \left(\frac{(1-z)(1-u)}{u} \right)^{\frac{1+\sqrt{1+24v}}{2}} B(z, u, v)$ and $z = 1 + t(1-u)/u$, the resulting differential equation is hypergeometric.

The differential equation (20) is exactly the same as the one for the number of passes in quickselect with median-of-three (see Proposition 1.3). The only difference between the expected number of passes in quickselect, as given in the work by Kirschenhofer et al. [17], and the number of ascendants in LBSTs relies on the initial conditions. The reason is that in the mentioned paper only one recursive call is counted if we want to select some element in a file of size ≤ 2 , while the average number of ascendants of the j^{th} node in a random LBST of size $n \leq 2$ is $3/2$ (for $j = 1$ and $j = 2$). Then $a_{n,j}$ and the expected number of passes to select the j^{th} element out of n differ in the constant term, exactly by $1/7$.

In a similar way, when differentiating the differential equation (19) two times w.r.t. v and setting $v = 1$, we get the differential equation for $\mathcal{A}_z^{(2)}(z, u)$, the generating function whose coefficients are the second factorial moments $a_{n,j}^{(2)}$ of the number of ascendants. Solving this differential equation and extracting the coefficients leads to the second factorial moments, which are given in the appendix.

In [27] the authors considered the expectation and variance of A_n in random LBSTs. To be more precise, they stated the problem in terms of unsuccessful search costs. Here, we are able to reproduce their results and extend them to higher order moments. Since we deal with ascendants of internal nodes, our results can be naturally stated in terms of successful search costs, and then translated to unsuccessful costs using Proposition 1.2. where $a_n = \mathbb{E}[A_n]$ is the expected number of ascendants of a random node in a random tree with n nodes.

Theorem 7.2. *Let*

$$\bar{A}_z(z, v) = \frac{\partial}{\partial z} \sum_{n,m} \mathbb{P}[A_n = m] z^n v^m.$$

Then

$$\begin{aligned} \bar{A}_z(z, v) &= \frac{v}{(1 - 2v)(1 - z)^2} \\ &\quad - \frac{v^2}{(1 - 2v)\Delta} \left((\Delta + 4v + 3)(1 - z)^{-(\Delta-1)/2} + (\Delta - 4v - 3)(1 - z)^{(\Delta+1)/2} \right), \end{aligned}$$

where $\Delta = \sqrt{1 + 48v}$.

Proof. The differential equation to be solved (from Equation (19), plugging $u = 1$) is

$$\frac{1}{6} \frac{\partial^2 \bar{A}_z}{\partial z^2} = \frac{2v}{(1 - z)^2} \bar{A}_z + \frac{v}{(1 - z)^4},$$

where $\bar{A}_z(z, v) = A_z(z, 1, v)$ and the initial conditions are $\bar{A}_z(0, v) = v$ and $\frac{\partial}{\partial z} \bar{A}_z(0, v) = 2v(1 + v)$. Recall that $\bar{A}_z(z, v) = A_z(z, 1, v)$. The solution of the differential equation above is the explicit form given in the theorem. \square

Extracting coefficients in exact form from there is quite difficult. However, as Philippe Flajolet kindly pointed to us, asymptotic information and most notably, the limiting probability distribution can be established [8, 15]. In this case, it follows that A_n converges in distribution (converges in law) to a Gaussian distribution, i.e.

$$\mathbb{P} \left[\frac{A_n - \frac{12}{7} \log n}{\sqrt{\frac{300}{343} \log n}} < x \right] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt + \mathcal{O} \left(\frac{1}{\sqrt{\log n}} \right).$$

This result follows from the asymptotic estimation for the average and the variance of A_n and the fact that $\bar{A}_z(z, v)$ is essentially a quasi-power of $[z^n] \bar{A}_z(z, v)$ in a neighborhood of $v = 1$, i.e.

$$[z^n] \bar{A}_z(z, v) = c(v) \cdot n^{(\Delta-3)/2} (1 + \mathcal{O}(1/\sqrt{n})),$$

and the error term is uniformly bounded. Using the expansion [6]

$$[z^n](1-z)^\alpha = \frac{n^{-\alpha-1}}{\Gamma(-\alpha)} \left(1 + \frac{\alpha(\alpha+1)}{2n} + \mathcal{O}\left(\frac{1}{n^2}\right) \right)$$

we get uniformly in the circle $|v-1| < \frac{1}{4}$

$$\begin{aligned} [z^n]\bar{A}_z(z, v) &= [z^n] - \frac{v^2}{(1-2v)\Delta} (\Delta + 4v + 3)(1-z)^{-(\Delta-1)/2} + \mathcal{O}(n) \\ &= -\frac{v^2(\Delta + 4v + 3)}{(1-2v)\Delta\Gamma(\frac{\Delta-1}{2})} \cdot n^{\frac{\Delta-3}{2}} \left(1 + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \right). \end{aligned}$$

Applying the following *quasi-power* theorem of Hwang [15, 7] leads immediately to the above given result.

Theorem 7.3. (Quasi-power theorem [H.-K. Hwang]) Assume that the Laplace transforms $\lambda_n(s) = \mathbb{E}[e^{sX_n}]$ of a sequence of random variables X_n are analytic in a disc $|s| < \rho$, for some $\rho > 0$, and satisfy there an expansion of the form

$$\lambda_n(s) = e^{\beta_n U(s) + V(s)} \left(1 + \mathcal{O}\left(\frac{1}{\kappa_n}\right) \right),$$

with $\beta_n, \kappa_n \rightarrow +\infty$, and $U(s), V(s)$ analytic in $|s| \leq \rho$. Assume also the variability condition,

$$U''(0) \neq 0.$$

Under these assumptions, the mean and variance of X_n satisfy

$$\mathbb{E}[X_n] = \beta_n U'(0) + V'(0) + \mathcal{O}(\kappa_n^{-1}), \quad \mathbb{V}[X_n] = \beta_n U''(0) + V''(0) + \mathcal{O}(\kappa_n^{-1}).$$

The distribution of X_n is asymptotically Gaussian and the speed of convergence to the Gaussian limit is $\mathcal{O}(\kappa_n^{-1} + \beta_n^{-1/2})$:

$$\mathbb{P}\left[\frac{X_n - \beta_n U'(0)}{\sqrt{\beta_n U''(0)}} \leq x \right] = \Phi(x) + \mathcal{O}\left(\frac{1}{\kappa_n} + \frac{1}{\sqrt{\beta_n}} \right).$$

$\Phi(x)$ denotes here the distribution function of the Gaussian normal distribution. □

The next step in our programme is to differentiate \bar{A}_z as many times as needed w.r.t. v and set $v = 1$, in order to get the generating functions for factorial moments.

Theorem 7.4. The expected number of ascendants $a_n = \mathbb{E}[A_n]$ of a random node in a random LBST of size n , when $n \geq 6$, is

$$a_n = \frac{12}{7} \left(1 + \frac{1}{n} \right) H_n - \frac{1}{49} \left(124 - \frac{9}{n} \right).$$

Proof. Let, as usual,

$$\bar{\mathcal{A}}^{(s)}(z) = \left. \frac{\partial^s \bar{A}_z}{\partial v^s} \right|_{v=1}.$$

To avoid cluttering the notation, we also let $\bar{\mathcal{A}}_z(z) = \bar{\mathcal{A}}^{(1)}(z)$. Here is the generating function for the expectations

$$\bar{\mathcal{A}}_z(z) = \frac{24}{7} \frac{1}{(1-z)^3} \log \frac{1}{1-z} + \frac{4}{49} (1-z)^{-3} + (1-z)^{-2} - \frac{4}{49} (1-z)^4.$$

Then we extract the $(n-1)^{\text{th}}$ coefficient and divide by n^2 to get the expected value of A_n ; recall that since we are averaging w.r.t. j and \bar{A}_z is already a partial derivative w.r.t. z , we have in fact

$$\mathbb{E}[A_n] = \frac{1}{n^2} [z^{n-1}] \bar{\mathcal{A}}_z(z).$$

□

Theorem 7.5. *The variance of the number of ascendants, A_n , of a random node in a random LBST with n nodes, or equivalently, the variance of the successful search cost for a random element in a LBST of size n is, when $n \geq 6$,*

$$\begin{aligned} \mathbb{V}[A_n] &= \frac{1}{343} \left(300 + \frac{2100}{n} - \frac{216}{n^2} \right) H_n \\ &\quad - \frac{144}{49} \left(1 + \frac{1}{n} \right) \left(\frac{H_n^2}{n} + H_n^{(2)} \right) + \frac{1}{2401} \left(10758 + \frac{2431}{n} - \frac{81}{n^2} \right) + \frac{2304}{343n n^6}. \end{aligned}$$

Proof. Analogously to what we did in the proof of the previous theorem, we compute the second derivative of $\bar{A}_z(z, v)$, and let $v = 1$. Then

$$\begin{aligned} \bar{\mathcal{A}}^{(2)}(z) &= \frac{288}{49} \frac{1}{(1-z)^3} \log^2 \frac{1}{1-z} + \left(-\frac{480}{343} \frac{1}{(1-z)^3} + \frac{96}{343} (1-z)^4 \right) \log \frac{1}{1-z} \\ &\quad + \frac{9988}{2401} \frac{1}{(1-z)^3} - 4 \frac{1}{(1-z)^2} - \frac{384}{2401} (1-z)^4. \end{aligned}$$

Extracting the coefficients is not as easy as before, but it is also doable, yielding the second factorial moment:

$$\begin{aligned} \mathbb{E}[A_n^2] &= \frac{144}{49} \left(1 + \frac{1}{n} \right) \left(H_n^2 - H_n^{(2)} \right) - \frac{1}{343} \left(3264 + \frac{1248}{n} \right) H_n \\ &\quad + \frac{1}{2401} \left(32210 - \frac{242}{n} \right) + \frac{2304}{343n n^6}. \end{aligned}$$

From here, the remaining computations are just mechanical. □

For higher order moments, i.e. $s > 2$, the procedure applies but the computations get messier. If we do only consider the main order term in $a_n^{(s)} = \mathbb{E}[A_n^s]$, then the result is much easier.

Theorem 7.6. *The s^{th} factorial moment of the number of ascendants, A_n , of a random node in a random LBST with n nodes, or equivalently, the s^{th} factorial moment of the successful search cost for a random element in a LBST of size n is, when $n \geq 6$,*

$$a_n^{(s)} = \left(\frac{12}{7} \right)^s \log^s n + \mathcal{O}(\log^{s-1} n).$$

ACKNOWLEDGEMENTS

We thank Philippe Flajolet for useful comments and suggestions. We also wish to thank the authors of the computer algebra system MAPLE who, although they might not know, greatly contributed to make this paper possible.

REFERENCES

- [1] G.M. Adel'son-Vel'skii and E.M. Landis. An algorithm for the organization of information. *Dokladi Akademii Nauk SSSR*, 146(2):263–266, 1962. English translation in *Soviet Math. Doklady* 3, 1962, 1259–1263.
- [2] C.R. Aragon and R.G. Seidel. Randomized search trees. In *Proc. of the 30th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 540–545, 1989.
- [3] S.R. Arora and W.T. Dent. Randomized binary search technique. *Comm. ACM*, 12(2):77–80, 1969.
- [4] C.J. Bell. *An Investigation into the Principles of the Classification and Analysis of Data on an Automatic Digital Computer*. PhD thesis, Leeds University, 1965.
- [5] G.G. Brown and B.O. Shubert. On random binary trees. *Mathematics of Operations Research*, 9(1):43–65, 1984.
- [6] Ph. Flajolet and A.M. Odlyzko. Singularity analysis of generating functions. *SIAM Journal on Discrete Mathematics*, 3(2):216–240, May 1990.
- [7] Ph. Flajolet and R. Sedgewick. The Average Case Analysis of Algorithms: Multivariate Asymptotics and Limit Distributions. *Rapport de recherche de l'INRIA #3162*, 1997.

- [8] Ph. Flajolet and M. Soria. General combinatorial schemas: Gaussian limit distributions and exponential tails. *Discrete Mathematics*, 114, 1993.
- [9] G.H. Gonnet and R. Baeza-Yates. *Handbook of Algorithms and Data Structures - In Pascal and C*. Addison-Wesley, 2nd edition, 1991.
- [10] R.L. Graham, D.E. Knuth, and O. Patashnik. *Concrete Mathematics*. Addison-Wesley, 1989.
- [11] L.J. Guibas and R. Sedgwick. A dichromatic framework for balanced trees. In *Proc. of the 19th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 8–21, 1978.
- [12] C.A.R. Hoare. Find (Algorithm 65). *Comm. ACM*, 4:321–322, 1961.
- [13] C.A.R. Hoare. Quicksort. *Comput. J.*, 5:10–15, 1962.
- [14] M. Hoshi and Ph. Flajolet. Page usage in a quadtree index. *BIT*, 32(3):384–402, 1992.
- [15] H.-K. Hwang. *Théorèmes limites pour les structures combinatoires et les fonctions arithmétiques*. PhD thesis, Ecole Polytechnique, 1994.
- [16] E. Kamke. *Differentialgleichungen: Lösungsmethoden und Lösungen*. Teubner, Stuttgart, 1977.
- [17] P. Kirschenhofer, C. Martínez, and H. Prodinger. Analysis of Hoare’s Find Algorithm with Median-of-three partition. *Random Structures & Algorithms*, 10:143–156, 1997.
- [18] P. Kirschenhofer and H. Prodinger. Comparisons in Hoare’s Find algorithm. *Combinatorics, Probability and Computing*, 7:111–120, 1998.
- [19] D.E. Knuth. Mathematical analysis of algorithms. In *Proc. of the 1971 IFIP Congress*, pages 19–27, Amsterdam, 1972. North-Holland.
- [20] D.E. Knuth. *The Art of Computer Programming: Sorting and Searching*, volume 3. Addison-Wesley, 1973.
- [21] J. Lent. *Probabilistic analysis of some searching and sorting algorithms*. PhD thesis, George Washington University, 1996.
- [22] G. Louchard. Exact and asymptotic distributions in digital and binary search trees. *Theoretical Informatics and Applications*, 21(4):479–496, 1987.
- [23] W.C. Lynch. More combinatorial properties of certain trees. *Comput. J.*, 7:299–302, 1965.
- [24] H.M. Mahmoud. *Evolution of Random Search Trees*. Wiley Interscience, 1992.
- [25] A. Panholzer. *Untersuchungen zur durchschnittlichen Gestalt gewisser Baumfamilien. Mit besonderer Berücksichtigung von Anwendungen in der Informatik*. PhD thesis, Technische Universität Wien, 1997.
- [26] P.V. Poblete. The analysis of heuristics for search trees. *Acta Informatica*, 30:233–248, 1993.
- [27] P.V. Poblete and J.I. Munro. The analysis of a fringe heuristic for binary search trees. *J. Algorithms*, 6:336–350, 1985.
- [28] S. Roura and C. Martínez. Randomization of search trees by subtree size. In J. Díaz and M. Serna, editors, *Proc. of the 4th European Symposium on Algorithms (ESA)*, volume 1136 of *LNCS*, pages 91–106. Springer, 1996.
- [29] B. Salvy and P. Zimmermann. Gfun: a Maple package for the manipulation of generating and holonomic functions in one variable. *ACM Transactions on Mathematical Software*, 20(2):163–177, 1994.
- [30] R. Sedgwick. *Quicksort*. Garland, New York, 1978.
- [31] R. Sedgwick. *Algorithms in C*. Addison-Wesley, 3rd edition, 1997.
- [32] R. Sedgwick and Ph. Flajolet. *An Introduction to the Analysis of Algorithms*. Addison-Wesley, 1996.
- [33] J.S. Vitter and Ph. Flajolet. Average-case analysis of algorithms and data structures. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, chapter 9. North-Holland, 1990.
- [34] A. Walker and D. Wood. Locally balanced binary trees. *Comput. J.*, 19(4):322–325, 1976.
- [35] H. Wilf. *Generatingfunctionology*. Academic Press, 1990.