

Self-Adjusting Binary Search Trees

DANIEL DOMINIC SLEATOR AND ROBERT ENDRE TARJAN

AT&T Bell Laboratories, Murray Hill, NJ

Abstract. The *splay* tree, a self-adjusting form of binary search tree, is developed and analyzed. The binary search tree is a data structure for representing tables and lists so that accessing, inserting, and deleting items is easy. On an n -node splay tree, all the standard search tree operations have an amortized time bound of $O(\log n)$ per operation, where by "amortized time" is meant the time per operation averaged over a worst-case sequence of operations. Thus splay trees are as efficient as balanced trees when total running time is the measure of interest. In addition, for sufficiently long access sequences, splay trees are as efficient, to within a constant factor, as static optimum search trees. The efficiency of splay trees comes not from an explicit structural constraint, as with balanced trees, but from applying a simple restructuring heuristic, called *splaying*, whenever the tree is accessed. Extensions of splaying give simplified forms of two other data structures: lexicographic or multidimensional search trees and link/cut trees.

Categories and Subject Descriptors: E.1 [Data]: Data Structures—*trees*; F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems—*sorting and searching*

General Terms: Algorithms, Theory

Additional Key Words and Phrases: Amortized complexity, balanced trees, multidimensional searching, network optimization, self-organizing data structures

1. Introduction

In this paper we apply the related concepts of *amortized complexity* and *self-adjustment* to binary search trees. We are motivated by the observation that the known kinds of efficient search trees have various drawbacks. Balanced trees, such as height-balanced trees [2, 22], weight-balanced trees [26], and B-trees [6] and their variants [5, 18, 19, 24] have a worst-case time bound of $O(\log n)$ per operation on an n -node tree. However, balanced trees are not as efficient as possible if the access pattern is nonuniform, and they also need extra space for storage of balance information. Optimum search trees [16, 20, 22] guarantee minimum average access time, but only under the assumption of fixed, known access probabilities and no correlation among accesses. Their insertion and deletion costs are also very high. Biased search trees [7, 8, 13] combine the fast average access time of optimum trees with the fast updating of balanced trees but have structural constraints even more complicated and harder to maintain than the constraints of balanced trees. Finger search trees [11, 14, 19, 23, 24] allow fast access in the vicinity of one or more "fingers" but require the storage of extra pointers in each node.

Authors' address: AT&T Bell Laboratories, 600 Mountain Avenue, Murray Hill, NJ 07974.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1985 ACM 0004-5411/85/0700-0652 \$00.75

These data structures are all designed to reduce the worst-case time per operation. However, in typical applications of search trees, not one but a sequence of operations is performed, and what matters is the total time the sequence takes, not the individual times of the operations. In such applications, a better goal is to reduce the amortized times of the operations, where by “amortized time” we mean the average time of an operation in a worst-case sequence of operations.

One way to obtain amortized efficiency is to use a “self-adjusting” data structure. We allow the structure to be in an arbitrary state, but during each operation we apply a simple restructuring rule intended to improve the efficiency of future operations. Examples of restructuring heuristics that give amortized efficiency are the move-to-front rule on linear lists [9, 30] and path compression in balanced trees [33, 37].

Self-adjusting data structures have several possible advantages over balanced or otherwise explicitly constrained structures:

- (i) In an amortized sense, ignoring constant factors, they are never much worse than constrained structures, and since they adjust according to usage, they can be much more efficient if the usage pattern is skewed.
- (ii) They need less space, since no balance or other constraint information is stored.
- (iii) Their access and update algorithms are conceptually simple and easy to implement.

Self-adjusting structures have two potential disadvantages:

- (i) They require more local adjustments, especially during accesses (look-up operations). (Explicitly constrained structures need adjusting only during updates, not during accesses.)
- (ii) Individual operations within a sequence can be expensive, which may be a drawback in real-time applications.

In this paper we develop and analyze the *splay* tree, a self-adjusting form of binary search tree. The restructuring heuristic used in splay trees is *splaying*, which moves a specified node to the root of the tree by performing a sequence of rotations along the (original) path from the node to the root. In an amortized sense and ignoring constant factors, splay trees are as efficient as both dynamically balanced trees and static optimum trees, and they may have even stronger optimality properties. In particular, we conjecture that splay trees are as efficient on any sufficiently long sequence of accesses as any form of dynamically updated binary search tree, even one tailored to the exact access sequence.

The paper contains seven sections. In Section 2 we define splaying and analyze its amortized complexity. In Section 3 we discuss update operations on splay trees. In Section 4 we study the practical efficiency and ease of implementation of splaying and some of its variants. In Section 5 we explore ways of reducing the amount of restructuring needed in splay trees. In Section 6 we use extensions of splaying to simplify two more complicated data structures: lexicographic or multidimensional search trees [15, 25] and link/cut trees [29, 34]. In Section 7 we make some final remarks and mention several open problems, including a formal version of our dynamic optimality conjecture. The appendix contains our tree terminology.

The work described here is a continuation of our research on amortized complexity and self-adjusting data structures, which has included an amortized analysis

of the move-to-front list update rule [9, 30] and the development of a self-adjusting form of heap [31]. Some of our results on self-adjusting heaps and search trees appeared in preliminary form in a conference paper [28]. A survey paper by the second author examines a variety of amortized complexity results [35].

2. Splay Trees

We introduce splay trees by way of a specific application. Consider the problem of performing a sequence of access operations on a set of items selected from a totally ordered universe. Each item may have some associated information. The input to each operation is an item; the output of the operation is an indication of whether the item is in the set, along with the associated information if it is. One way to solve this problem is to represent the set by a *binary search tree*. This is a binary tree containing the items of the set, one item per node, with the items arranged in *symmetric order*: If x is a node containing an item i , the left subtree of x contains only items less than i and the right subtree only items greater than i . The *symmetric-order position* of an item is one plus the number of items preceding it in symmetric order in the tree.

The “search” in “binary search tree” refers to the ability to access any item in the tree by searching down from the root, branching left or right at each step according to whether the item to be found is less than or greater than the item in the current node, and stopping when the node containing the item is reached. Such a search takes $\Theta(d)$ time, where d is the depth of the node containing the accessed item.

If accessing items is the only operation required, then there are better solutions than binary search trees, e.g., hashing. However, as we shall see in Section 3, binary search trees also support several useful update operations. Furthermore, we can extend binary search trees to support accesses by symmetric-order position. To do this, we store in each node the number of descendants of the node. Alternatively, we can store in each node the number of nodes in its left subtree and store in a tree header the number of nodes in the entire tree.

When referring to a binary search tree formally, as in a computer program, we shall generally denote the tree by a pointer to its root; a pointer to the null node denotes an empty tree. When analyzing operations on binary search trees, we shall use n to denote the number of nodes and m to denote the total number of operations.

Suppose we wish to carry out a sequence of access operations on a binary search tree. For the total access time to be small, frequently accessed items should be near the root of the tree often. Our goal is to devise a simple way of restructuring the tree after each access that moves the accessed item closer to the root, on the plausible assumption that this item is likely to be accessed again soon. As an $O(1)$ -time restructuring primitive, we can use *rotation*, which preserves the symmetric order of the tree. (See Figure 1.)

Allen and Munro [4] and Bitner [10] proposed two restructuring heuristics (see Figure 2):

Single rotation. After accessing an item i in a node x , rotate the edge joining x to its parent (unless x is the root).

Move to root. After accessing an item i in a node x , rotate the edge joining x to its parent, and repeat this step until x is the root.

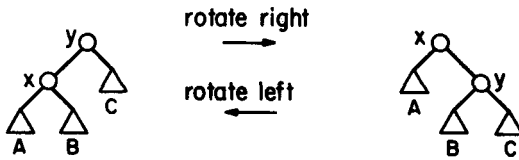


FIG. 1. Rotation of the edge joining nodes x and y . Triangles denote subtrees. The tree shown can be a subtree of a larger tree.

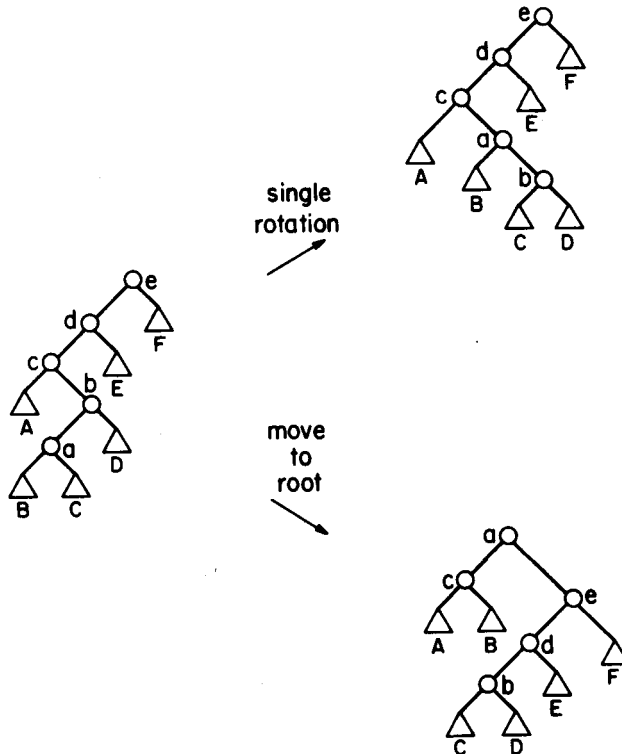


FIG. 2. Restructuring heuristics. The node accessed is a .

Unfortunately, neither of these heuristics is efficient in an amortized sense: for each, there are arbitrarily long access sequences such that the time per access is $O(n)$ [4]. Allen and Munro did show that the move-to-root heuristic has an asymptotic average access time that is within a constant factor of minimum, but only under the assumption that the access probabilities of the various items are fixed and the accesses are independent. We seek heuristics that have much stronger properties.

Our restructuring heuristic, called *splaying*, is similar to move-to-root in that it does rotations bottom-up along the access path and moves the accessed item all the way to the root. But it differs in that it does the rotations in pairs, in an order that depends on the structure of the access path. To splay a tree at a node x , we repeat the following *splaying step* until x is the root of the tree (see Figure 3):

Splaying Step

Case 1 (zig). If $p(x)$, the parent of x , is the tree root, rotate the edge joining x with $p(x)$. (This case is terminal.)

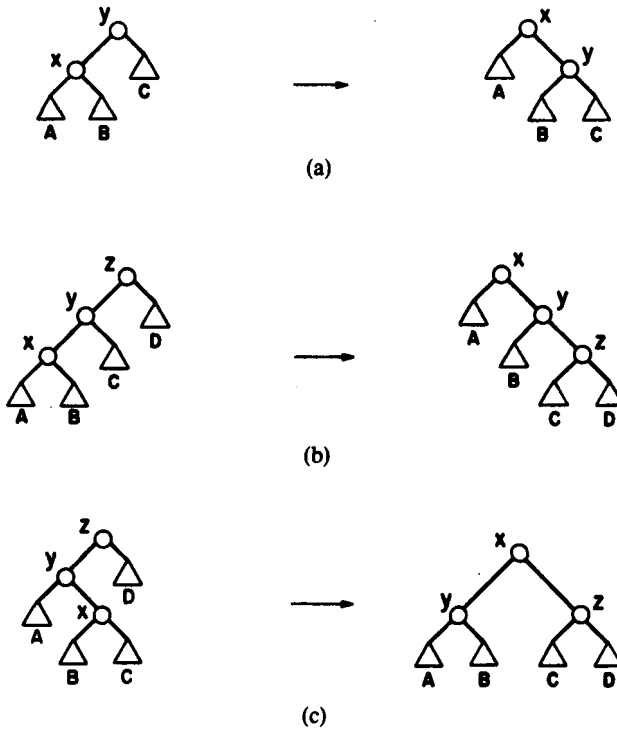


FIG. 3. A splaying step. The node accessed is x . Each case has a symmetric variant (not shown). (a) Zig: terminating single rotation. (b) Zig-zig: two single rotations. (c) Zig-zag: double rotation.

Case 2 (zig-zig). If $p(x)$ is not the root and x and $p(x)$ are both left or both right children, rotate the edge joining $p(x)$ with its grandparent $g(x)$ and then rotate the edge joining x with $p(x)$.

Case 3 (zig-zag). If $p(x)$ is not the root and x is a left child and $p(x)$ a right child, or vice-versa, rotate the edge joining x with $p(x)$ and then rotate the edge joining x with the new $p(x)$.

Splaying at a node x of depth d takes $\Theta(d)$ time, that is, time proportional to the time to access the item in x . Splaying not only moves x to the root, but roughly halves the depth of every node along the access path. (See Figures 4 and 5.) This halving effect makes splaying efficient and is a property not shared by other, simpler heuristics, such as move to root. Producing this effect seems to require dealing with the zig-zig and zig-zag cases differently.

We shall analyze the amortized complexity of splaying by using a *potential function* [31, 35] to carry out the amortization. The idea is to assign to each possible configuration of the data structure a real number called its *potential* and to define the *amortized time* a of an operation by $a = t + \Phi' - \Phi$, where t is the actual time of the operation, Φ is the potential before the operation, and Φ' is the potential after the operation. With this definition, we can estimate the total time of a sequence of m operations by

$$\sum_{j=1}^m t_j = \sum_{j=1}^m (a_j + \Phi_{j-1} - \Phi_j) = \sum_{j=1}^m a_j + \Phi_0 - \Phi_m,$$

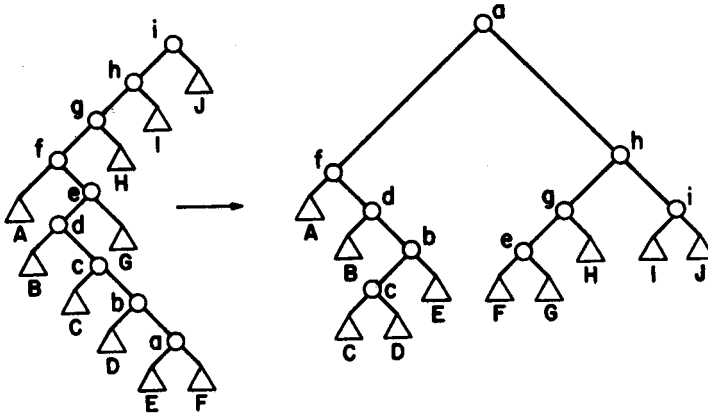


FIG. 4. Splaying at node *a*.

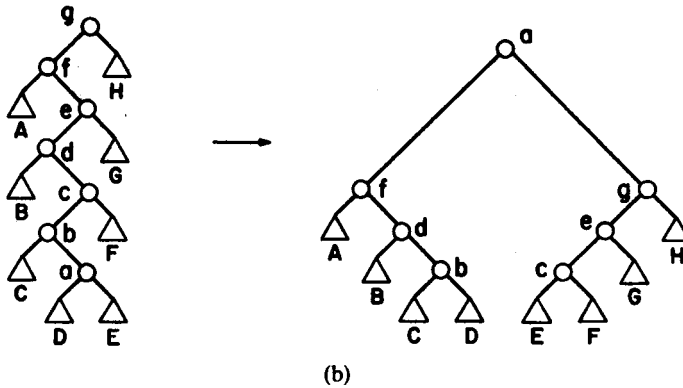
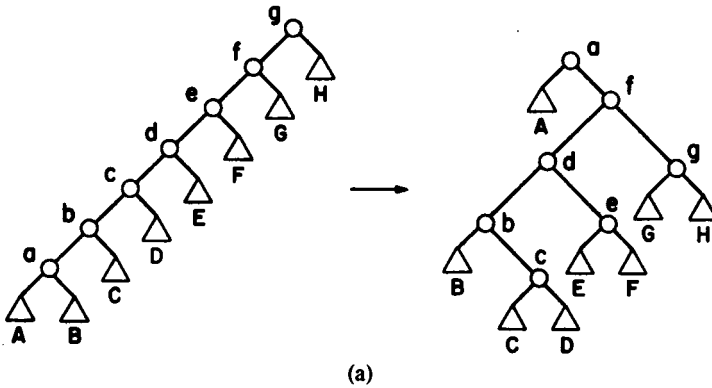


FIG. 5. Extreme cases of splaying. (a) All zig-zig steps. (b) All zig-zag steps.

where t_j and a_j are the actual and amortized times of operation j , Φ_0 is the initial potential, and Φ_j for $j \geq 1$ is the potential after operation j . That is, the total actual time equals the total amortized time plus the net decrease in potential from the initial to the final configuration. If the final potential is no less than the initial potential, then the total amortized time is an upper bound on the total actual time.

To define the potential of a splay tree, we assume that each item i has a positive weight $w(i)$, whose value is arbitrary but fixed. We define the size $s(x)$ of a node x in the tree to be the sum of the individual weights of all items in the subtree rooted at x . We define the rank $r(x)$ of node x to be $\log s(x)$.¹ Finally, we define the potential of the tree to be the sum of the ranks of all its nodes. As a measure of the running time of a splaying operation, we use the number of rotations done, unless there are no rotations, in which case we charge one for the splaying.

LEMMA 1 (ACCESS LEMMA). *The amortized time to splay a tree with root t at a node x is at most $3(r(t) - r(x)) + 1 = O(\log(s(t)/s(x)))$.*

PROOF. If there are no rotations, the bound is immediate. Thus suppose there is at least one rotation. Consider any splaying step. Let s and s' , r and r' denote the size and rank functions just before and just after the step, respectively. We show that the amortized time for the step is at most $3(r'(x) - r(x)) + 1$ in case 1 and at most $3(r'(x) - r(x))$ in case 2 or case 3. Let y be the parent of x and z be the parent of y (if it exists) before the step.

Case 1. One rotation is done, so the amortized time of the step is

$$\begin{aligned} 1 + r'(x) + r'(y) - r(x) - r(y) & \quad \text{since only } x \text{ and } y \\ & \quad \text{can change rank} \\ \leq 1 + r'(x) - r(x) & \quad \text{since } r(y) \geq r'(y) \\ \leq 1 + 3(r'(x) - r(x)) & \quad \text{since } r'(x) \geq r(x). \end{aligned}$$

Case 2. Two rotations are done, so the amortized time of the step is

$$\begin{aligned} 2 + r'(x) + r'(y) + r'(z) \\ - r(x) - r(y) - r(z) & \quad \text{since only } x, y, \text{ and } z \\ & \quad \text{can change rank} \\ = 2 + r'(y) + r'(z) - r(x) - r(y) & \quad \text{since } r'(x) = r(z) \\ \leq 2 + r'(x) + r'(z) - 2r(x) & \quad \text{since } r'(x) \geq r'(y) \\ & \quad \text{and } r(y) \geq r(x). \end{aligned}$$

We claim that this last sum is at most $3(r'(x) - r(x))$, that is, that $2r'(x) - r(x) - r'(z) \geq 2$. The convexity of the log function implies that $\log x + \log y$ for $x, y > 0$, $x + y \leq 1$ is maximized at value -2 when $x = y = \frac{1}{2}$. It follows that $r(x) + r'(z) - 2r'(x) = \log(s(x)/s'(x)) + \log(s'(z)/s'(x)) \leq -2$, since $s(x) + s'(z) \leq s'(x)$. Thus the claim is true, and the amortized time in case 2 is at most $3(r'(x) - r(x))$.

Case 3. The amortized time of the step is

$$\begin{aligned} 2 + r'(x) + r'(y) + r'(z) \\ - r(x) - r(y) - r(z) & \quad \text{since } r'(x) = r(z) \\ \leq 2 + r'(y) + r'(z) - 2r(x) & \quad \text{and } r(x) \leq r(y). \end{aligned}$$

We claim that this last sum is at most $2(r'(x) - r(x))$, that is, that $2r'(x) - r'(y) - r'(z) \geq 2$. This follows as in case 2 from the inequality $s'(y) + s'(z) \leq s'(x)$. Thus the amortized time in case 3 is at most $2(r'(x) - r(x)) \leq 3(r'(x) - r(x))$.

The lemma follows by summing the amortized time estimates for all the splaying steps, since the sum telescopes to yield an estimate of at most $3(r'(x) - r(x)) + 1 = 3(r(t) - r(x)) + 1$, where r and r' are the rank functions before and after the entire splaying operation, respectively. \square

¹ Throughout this paper we use binary logarithms.

The analysis in Lemma 1 shows that the zig-zig case is the expensive case of a splaying step. The analysis differs from our original analysis [28] in that the definition of rank uses the continuous instead of the discrete logarithm. This gives us a bound that is tighter by a factor of two. The idea of tightening the analysis in this way was also discovered independently by Huddleston [17].

The weights of the items are parameters of the analysis, not of the algorithm: Lemma 1 holds for *any* assignment of positive weights to items. By choosing weights cleverly, we can obtain surprisingly strong results from Lemma 1. We shall give four examples. Consider a sequence of m accesses on an n -node splay tree. In analyzing the running time of such a sequence, it is useful to note that if the weights of all items remain fixed, then the net decrease in potential over the sequence is at most $\sum_{i=1}^n \log(W/w(i))$, where $W = \sum_{i=1}^n w(i)$, since the size of the node containing item i is at most W and at least $w(i)$.

THEOREM 1 (BALANCE THEOREM). *The total access time is $O((m + n) \log n + m)$.*

PROOF. Assign a weight of $1/n$ to each item. Then $W = 1$, the amortized access time is at most $3 \log n + 1$ for any item, and the net potential drop over the sequence is at most $n \log n$. The theorem follows. \square

For any item i , let $q(i)$ be the access frequency of item i , that is, the total number of times i is accessed.

THEOREM 2 (STATIC OPTIMALITY THEOREM). *If every item is accessed at least once, then the total access time is*

$$O\left(m + \sum_{i=1}^n q(i) \log\left(\frac{m}{q(i)}\right)\right).$$

PROOF. Assign a weight of $q(i)/m$ to item i . Then $W = 1$, the amortized access time of item i is $O(\log(m/q(i)))$, and the net potential drop over the sequence is at most $\sum_{i=1}^n \log(m/q(i))$. The theorem follows. \square

Assume that the items are numbered from 1 through n in symmetric order. Let the sequence of accessed items be i_1, i_2, \dots, i_m .

THEOREM 3 (STATIC FINGER THEOREM). *If f is any fixed item, the total access time is $O(n \log n + m + \sum_{j=1}^m \log(|i_j - f| + 1))$.*

PROOF. Assign a weight of $1/(|i - f| + 1)^2$ to item i . Then $W \leq 2 \sum_{k=1}^{\infty} 1/k^2 = O(1)$, the amortized time of the j th access is $O(\log(|i_j - f| + 1))$, and the net potential drop over the sequence is $O(n \log n)$, since the weight of any item is at least $1/n^2$. The theorem follows. \square

We can obtain another interesting result by changing the item weights as the accesses take place. Number the accesses from 1 to m in the order they occur. For any access j , let $t(j)$ be the number of different items accessed before access j since the last access of item i_j , or since the beginning of the sequence if j is the first of item i_j . (Note that $t(j) + i$ is the position of item i_j in a linear list maintained by the move-to-front heuristic [30] and initialized in order of first access.)

THEOREM 4 (WORKING SET THEOREM). *The total access time is $O(n \log n + m + \sum_{j=1}^m \log(t(j) + 1))$.*

PROOF. Assign the weights $1, 1/4, 1/9, \dots, 1/n^2$ to the items in order by first access. (The item accessed earliest gets the largest weight; any items never accessed get the smallest weights.) As each access occurs, redefine the weights as follows. Suppose the weight of item i_j during access j is $1/k^2$. After access j , assign weight 1 to i_j , and for each item i having a weight of $1/(k')^2$ with $k' < k$, assign weight $1/(k' + 1)^2$ to i . This reassignment permutes the weights $1, 1/4, 1/9, \dots, 1/n^2$ among the items. Furthermore, it guarantees that, during access j , the weight of item i_j will be $1/(t(j) + 1)^2$. We have $W = \sum_{k=1}^n 1/k^2 = O(1)$, so the amortized time of access j is $O(\log(t(j) + 1))$. The weight reassignment after an access increases the weight of the item in the root (because splaying at the accessed item moves it to the root) and only decreases the weights of other items in the tree. The size of the root is unchanged, but the sizes of other nodes can decrease. Thus the weight reassignment can only decrease the potential, and the amortized time for weight reassignment is either zero or negative. The net potential drop over the sequence is $O(n \log n)$. The theorem follows. \square

Let us interpret the meaning of these theorems. The balance theorem states that on a sufficiently long sequence of accesses a splay tree is as efficient as any form of uniformly balanced tree. The static optimality theorem implies that a splay tree is as efficient as any fixed search tree, including the optimum tree for the given access sequence, since by a standard theorem of information theory [1] the total access time for any fixed tree is $\Omega(m + \sum_{i=1}^n q(i) \log(m/q(i)))$. The static finger theorem states that splay trees support accesses in the vicinity of a fixed finger with the same efficiency as finger search trees. The working set theorem states that the time to access an item can be estimated as the logarithm of one plus the number of different items accessed since the given item was last accessed. That is, the most recently accessed items, which can be thought of as forming the "working set," are the easiest to access. All these results are to within a constant factor.

Splay trees have all these behaviors automatically; the restructuring heuristic is blind to the properties of the access sequence and to the global structure of the tree. Indeed, splay trees have all these behaviors simultaneously; at the cost of a constant factor we can combine all four theorems into one.

THEOREM 5 (UNIFIED THEOREM). *The total time of a sequence of m accesses on an n -node splay tree is*

$$O\left(n \log n + m + \sum_{j=1}^m \log \min \left\{ \frac{m}{q(i_j)}, |i_j - f| + 1, t(j) + 1 \right\}\right),$$

where f is any fixed item.

PROOF. Assign to each item a weight equal to the sum of the weights assigned to it in Theorems 2–4 and combine the proofs of these theorems. \square

Remark. Since $|i_j - f| < n$, Theorem 5 implies Theorem 1 as well as Theorems 2–4. If each item is accessed at least once, the additive term $n \log n$ in the bound of Theorem 5 can be dropped.

3. Update Operations on Splay Trees

Using splaying, we can implement all the standard update operations on binary search trees. We consider the following operations:

access(i, t): If item i is in tree t , return a pointer to its location; otherwise, return a pointer to the null node.

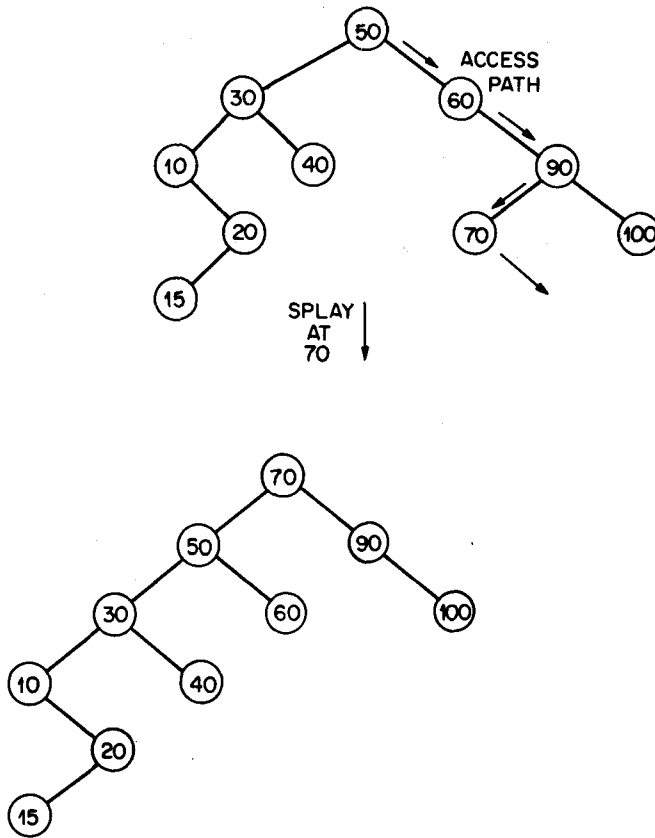


FIG. 6. Splaying after an unsuccessful access of item 80.

- insert*(i, t): Insert item i in tree t , assuming that it is not there already.
- delete*(i, t): Delete item i from tree t , assuming that it is present.
- join*(t_1, t_2): Combine trees t_1 and t_2 into a single tree containing all items from both trees and return the resulting tree. This operation assumes that all items in t_1 are less than all those in t_2 and destroys both t_1 and t_2 .
- split*(i, t): Construct and return two trees t_1 and t_2 , where t_1 contains all items in t less than or equal to i , and t_2 contains all items in t greater than i . This operation destroys t .

We can carry out these operations on splay trees as follows. To perform *access*(i, t), we search down from the root of t , looking for i . If the search reaches a node x containing i , we complete the access by splaying at x and returning a pointer to x . If the search reaches the null node, indicating that i is not in the tree, we complete the access by splaying at the last nonnull node reached during the search (the node from which the search ran off the bottom of the tree) and returning a pointer to null. If the tree is empty, we omit the splaying operation. (See Figure 6.)

Splaying's effect of moving a designated node to the root considerably simplifies the updating of splay trees. It is convenient to implement *insert* and *delete* using *join* and *split*. To carry out *join*(t_1, t_2), we begin by accessing the largest item, say

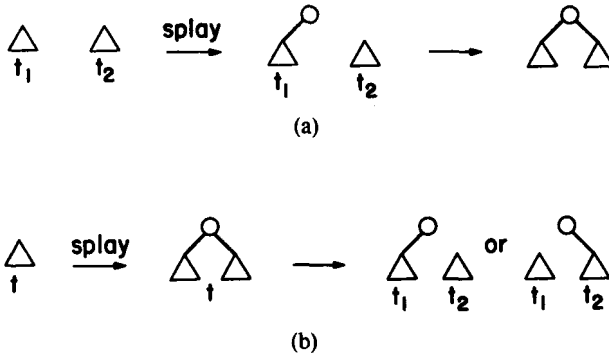


FIG. 7. Implementation of *join* and *split*: (a) *join*(t_1, t_2). (b) *split*(i, t).

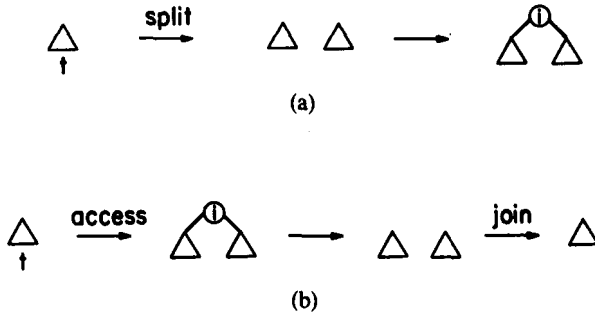


FIG. 8. Implementation of insertion and deletion using *join* and *split*: (a) *insert*(i, t). (b) *delete*(i, t).

i , in t_1 . After the access, the root of t_1 contains i and thus has a null right child. We complete the join by making t_2 the right subtree of this root and returning the resulting tree. To carry out *split*(i, t), we perform *access*(i, t) and then return the two trees formed by breaking either the left link or the right link from the new root of t , depending on whether the root contains an item greater than i or not greater than i . (See Figure 7.) In both *join* and *split* we must deal specially with the case of an empty input tree (or trees).

To carry out *insert*(i, t), we perform *split*(i, t) and then replace t by a tree consisting of a new root node containing i , whose left and right subtrees are the trees t_1 and t_2 returned by the split. To carry out *delete*(i, t), we perform *access*(i, t) and then replace t by the join of its left and right subtrees. (See Figure 8.)

There are alternative implementations of *insert* and *delete* that have slightly better amortized time bounds. To carry out *insert*(i, t), we search for i , then replace the pointer to null reached during the search by a pointer to a new node containing i , and finally splay the tree at the new node. To carry out *delete*(i, t), we search for the node containing i . Let this node be x and let its parent be y . We replace x as a child of y by the join of the left and right subtrees of x , and then we splay at y . (See Figure 9.)

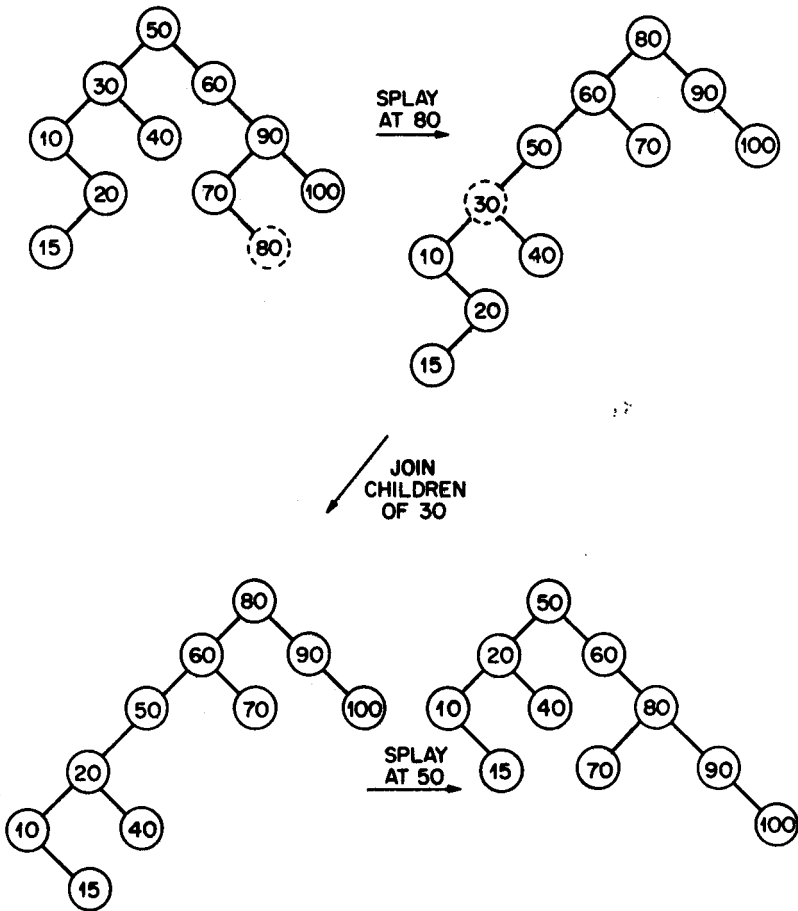


FIG. 9. Alternative implementations of insertion and deletion. Insertion of 80 followed by deletion of 30.

In analyzing the amortized complexity of the various operations on splay trees, let us assume that we begin with a collection of empty trees and that every item is only in one tree at a time. We define the potential of a collection of trees to be the sum of the potentials of the trees plus the sum of the logarithms of the weights of all items not currently in a tree. Thus the net potential drop over a sequence of operations is at most $\sum_{i \in U} \log(w(i)/w'(i))$, where U is the universe of possible items and w and w' , respectively, are the initial and final weight functions. In particular, if the item weights never change, the net potential change over the sequence is nonnegative.

For any item i in a tree t , let $i-$ and $i+$, respectively, denote the item preceding i and the item following i in t (in symmetric order). If $i-$ is undefined, let $w(i-) = \infty$; if $i+$ is undefined, let $w(i+) = \infty$.

LEMMA 2 (UPDATE LEMMA). *The amortized times of the splay tree operations have the following upper bounds, where W is the total weight of the items in the*

tree or trees involved in the operation:

$$access(i, t): \begin{cases} 3 \log\left(\frac{W}{w(i)}\right) + 1 & \text{if } i \text{ is in } t; \\ 3 \log\left(\frac{W}{\min\{w(i-), w(i+)\}}\right) + 1 & \text{if } i \text{ is not in } t. \end{cases}$$

$$join(t_1, t_2): 3 \log\left(\frac{W}{w(i)}\right) + O(1), \quad \text{where } i \text{ is the last item in } t_1.$$

$$split(i, t): \begin{cases} 3 \log\left(\frac{W}{w(i)}\right) + O(1) & \text{if } i \text{ is in } t; \\ 3 \log\left(\frac{W}{\min\{w(i-), w(i+)\}}\right) + O(1) & \text{if } i \text{ is not in } t. \end{cases}$$

$$insert(i, t): 3 \log\left(\frac{W - w(i)}{\min\{w(i-), w(i+)\}}\right) + \log\left(\frac{W}{w(i)}\right) + O(1).$$

$$delete(i, t): 3 \log\left(\frac{W}{w(i)}\right) + 3 \log\left(\frac{W - w(i)}{w(i-)}\right) + O(1).$$

Increasing the weight of the item in the root of a tree t by an amount δ takes at most $\log(1 + \delta/W)$ amortized time, and decreasing the weight of any item takes negative amortized time.

PROOF. These bounds follow from Lemma 1 and a straightforward analysis of the potential change caused by each operation. Let s be the size function just before the operation. In the case of an *access* or *split*, the amortized time is at most $3 \log(s(t)/s(x)) + 1$ by Lemma 1, where x is the node at which the tree is splayed. If item i is in the tree, it is in x , and $s(x) \geq w(i)$. If i is not in the tree, either $i-$ or $i+$ is in the subtree rooted at x , and $s(x) \geq \min\{w(i-), w(i+)\}$. This gives the bound on *access* and *split*. The bound on *join* is immediate from the bound on *access*: the splaying time is at most $3 \log(s(t_1)/w(i)) + 1$, and the increase in potential caused by linking t_1 and t_2 is

$$\log\left(\frac{s(t_1) + s(t_2)}{s(t_1)}\right) \leq 3 \log\left(\frac{W}{s(t_1)}\right).$$

(We have $W = s(t_1) + s(t_2)$.) The bound on *insert* also follows from the bound on *access*; the increase in potential caused by adding a new root node containing i is

$$\log\left(\frac{s(t) + w(i)}{w(i)}\right) = \log\left(\frac{W}{w(i)}\right).$$

The bound on *delete* is immediate from the bounds on *access* and *join*. \square

Remark. The alternative implementations of insertion and deletion have the following upper bounds on their amortized times:

$$\text{insert}(i, t): \quad 2 \log\left(\frac{W - w(i)}{\min\{w(i-), w(i+)\}}\right) + \log\left(\frac{W}{w(i)}\right) + O(1).$$

$$\text{delete}(i, t): \quad 3 \log\left(\frac{W - w(i)}{\min\{w(i-), w(i)\}}\right) + O(1).$$

These bounds follow from a modification of the proof of Lemma 1. For the case of equal-weight items, the alternative forms of insertion and deletion each take at most $3 \log n + O(1)$ amortized time on a n -node tree. This bound was obtained independently by Huddleston [17] for the same insertion algorithm and a slightly different deletion algorithm.

The bounds in Lemma 2 are comparable to the bounds for the same operations on biased trees [7, 8, 13], but the biased tree algorithms depend explicitly on the weights. By using Lemma 2, we can extend the theorems of Section 2 in various ways to include update operations. (An exception is that we do not see how to include deletion in a theorem analogous to Theorem 3.) We give here only the simplest example of such an extension.

THEOREM 6 (BALANCE THEOREM WITH UPDATES). *A sequence of m arbitrary operations on a collection of initially empty splay trees takes $O(m + \sum_{j=1}^m \log n_j)$ time, where n_j is the number of items in the tree or trees involved in operation j .*

PROOF. Assign to each item a fixed weight of one and apply Lemma 2. \square

We can use splay trees as we would use any other binary search tree; for example, we can use them to implement various abstract data structures consisting of sorted sets or lists subject to certain operations. Such structures include dictionaries, which are sorted sets with access, insertion, and deletion, and *concatenable queues*, which are lists with access by position, insertion, deletion, concatenation, and splitting [3, 22]. We investigate two further applications of splaying in Section 6.

4. Implementations of Splaying and Its Variants

In this section we study the implementation of splaying and some of its variants. Our aim is to develop a version that is easy to program and efficient in practice. As a programming notation, we shall use a version of Dijkstra's guarded command language [12], augmented by the addition of procedures and "initializing guards" (G. Nelson, private communication). We restrict our attention to successful accesses, that is, accesses of items known to be in the tree.

Splaying, as we have defined it, occurs during a second, bottom-up pass over an access path. Such a pass requires the ability to get from any node on the access path to its parent. To make this possible, we can save the access path as it is traversed (either by storing it in an auxiliary stack or by using "pointer reversal" to encode it in the tree structure), or we can maintain parent pointers for every node in the tree. If space is expensive, we can obtain the effect of having parent pointers without using extra space, by storing in each node a pointer to its leftmost child and a pointer to its right sibling, or to its parent if it has no right sibling. (See Figure 10.) This takes only two pointers per node, the same as the standard left-child-right-child representation, and allows accessing the left child, right child, or

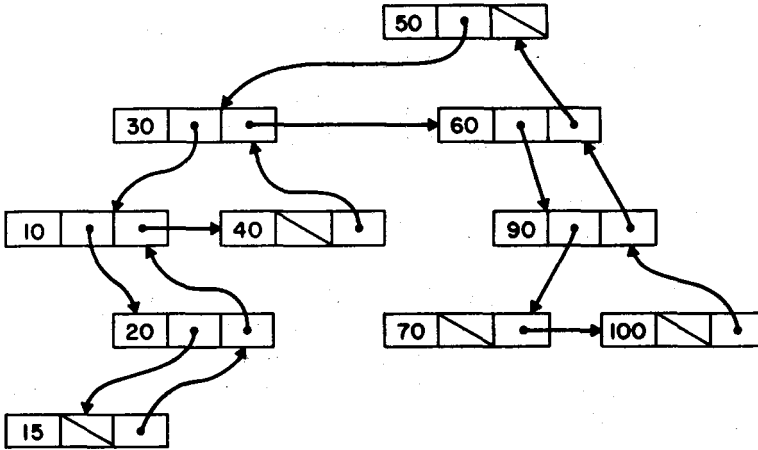


FIG. 10. Leftmost-child-right-sibling representation of the original tree in Figure 6.

parent of any node in at most two steps. The drawback of this representation is loss of time efficiency.

We shall assume that parent pointers are stored explicitly; however, our programs can easily be modified to use other ways of retrieving the access path. The following procedure is a straightforward implementation of splaying. To avoid an explicit test for $p(x) = \text{null}$, we assume that $p(\text{null}) = \text{left}(\text{null}) = \text{right}(\text{null}) = \text{null}$.

```

procedure splay( $x$ );
  {  $p(\text{null}) = \text{left}(\text{null}) = \text{right}(\text{null}) = \text{null}$  }
  do  $x = \text{left}(p(x)) \rightarrow$ 
    if  $g(x) = \text{null} \rightarrow \text{rotate right}(p(x))$ 
    |  $p(x) = \text{left}(g(x)) \rightarrow \text{rotate right}(g(x)); \text{rotate right}(p(x))$ 
    |  $p(x) = \text{right}(g(x)) \rightarrow \text{rotate right}(p(x)); \text{rotate left}(p(x))$ 
    fi
  |  $x = \text{right}(p(x)) \rightarrow$ 
    if  $g(x) = \text{null} \rightarrow \text{rotate left}(p(x))$ 
    |  $p(x) = \text{right}(g(x)) \rightarrow \text{rotate left}(g(x)); \text{rotate left}(p(x))$ 
    |  $p(x) = \text{left}(g(x)) \rightarrow \text{rotate left}(p(x)); \text{rotate right}(p(x))$ 
    fi
  od {  $p(x) = \text{null}$  }
end splay;

```

The grandparent function g is defined as follows:

```

function  $g(x)$ ;
   $g := p(p(x))$ 
end  $g$ ;

```

The procedure $\text{rotate left}(y)$ rotates the edge joining y and its right child. The procedure $\text{rotate right}(y)$, whose implementation we omit, is symmetric.

Remark. In the following procedure, the initializing guard " $x, z: x = \text{right}(y)$ and $z = p(y)$," which is always true, declares two local variables, x and z , and initializes them to $\text{right}(y)$ and $p(y)$, respectively.

```

procedure rotate left(y);
  if x, z: x = right(y) and z = p(y) →
    if z = null → skip
    | left(z) = y → left(z) := x
    | right(z) = y → right(z) := x
  fi;
  left(x), right(y) := y, left(x);
  p(x), p(y) := z, x;
  if right(y) = null → skip | right(y) ≠ null → p(right(y)) := y fi
fi
end rotate left;

```

Inspection of the splay program raises the possibility of simplifying it by omitting the second rotation in the zig-zag case. The resulting program, suggested by M. D. McIlroy (private communication), appears below.

```

procedure simple splay(x);
  {p(null) = left(null) = right(null) = null}
  do x = left(p(x)) →
    if p(x) = left(g(x)) → rotate right(g(x))
    | p(x) ≠ left(g(x)) → skip
  fi;
  rotate right(p(x))
  | x = right(p(x)) →
  if p(x) = right(g(x)) → rotate left(g(x))
  | p(x) ≠ right(g(x)) → skip
  fi;
  rotate left(p(x))
  od{p(x) = null}
end simple splay;

```

An amortized analysis of simple splaying shows that Lemma 1 holds with a constant factor of 3.16+ in place of three. Thus the method, though simpler than the original method, has a slightly higher bound on the number of rotations.

The advantage of implementing splaying using rotation as a primitive is that we can encapsulate all the restructuring, including any required updating of auxiliary fields, in the rotation procedures. The disadvantage is that many unnecessary pointer assignments are done. We can achieve greater efficiency by eliminating the superfluous assignments by expanding the rotation procedures in-line, simplifying, and keeping extra information in the control state.

A bottom-up splaying method is appropriate if we have direct access to the node at which the splaying is to occur. We shall see an example of this in Section 6. However, splaying as used in Sections 2 and 3 only occurs immediately after an access, and it is more efficient to splay on the way down the tree during the access. We present a top-down version of splaying that works in this way.

Suppose we wish to access an item i in a tree. During the access and concurrent splaying operation, the tree is broken into three parts: a *left tree*, a *middle tree*, and a *right tree*. The left and right trees contain all items in the original tree so far known to be less than i and greater than i , respectively. The middle tree consists of the subtree of the original tree rooted at the current node on the access path. Initially the current node is the tree root and the left and right trees are empty. To do the splaying, we search down from the root looking for i , two nodes at a time, breaking links along the access path and adding each subtree thus detached to the bottom right of the left tree or the bottom left of the right tree, with the proviso

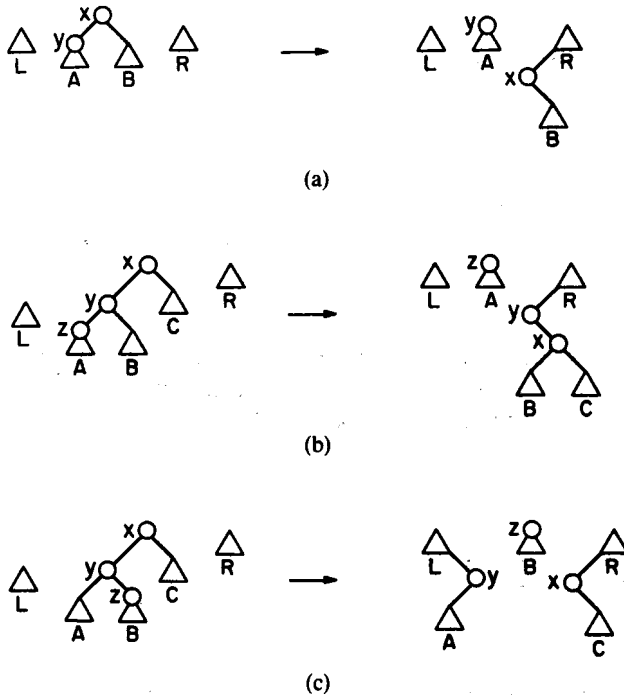


FIG. 11. Top-down splaying step. Symmetric variants of cases are omitted. Initial left tree is L , right tree is R . Labeled nodes are on the access path. (a) Zig: Node y contains the accessed item. (b) Zig-zig. (c) Zig-zag.

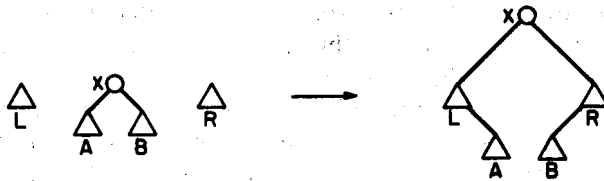


FIG. 12. Completion of top-down splaying. Node x contains the accessed item.

that if we take two steps down in the same direction (left-left or right-right) we do a rotation before breaking a link. More formally, the splaying operation consists of repeatedly applying the appropriate case among those shown in Figure 11 until none applies, then reassembling the three remaining trees as shown in Figure 12. A similar top-down restructuring method, but without the rotations and consequently with poor amortized efficiency, was proposed by Stephenson [32].

The following program implements this method. The program uses three local variables, t , l , and r , denoting the current vertex, the last node in symmetric order in the left tree, and the first node in symmetric order in the right tree, respectively. The updating of these variables occurs in the primitive procedures. These are *rotate left* and *rotate right*, which rotate the edge joining t and its right or left child, respectively; *link left* and *link right*, which break the link between t and its left or right child and attach the resulting tree to the left or right tree, respectively; and *assemble*, which completes the splaying by performing the assembly of Figure 12. The program contains an initializing guard that declares l and r and initializes

them both to null. After the first left link, the right child of null is the root of the left tree; after the first right link, the left child of null is the root of the right tree.

```

procedure top-down splay(i, t);
  if l, r: l = r = null  $\rightarrow$ 
    left(null) := right(null) := null;
  do i < item(t)  $\rightarrow$ 
    if i = item(left(t))  $\rightarrow$  link right
      | i < item(left(t))  $\rightarrow$  rotate right; link right
      | i > item(left(t))  $\rightarrow$  link right; link left
    fi
    | i > item(t)  $\rightarrow$ 
      if i = item(right(t))  $\rightarrow$  link left
        | i > item(right(t))  $\rightarrow$  rotate left; link left
        | i < item(right(t))  $\rightarrow$  link left; link right
      fi
    od{i = item(t)};
  assemble
fi
end top-down splay;

```

Implementations of *rotate left*, *link left*, and *assemble* appear below; *rotate right* and *link right*, whose definitions we omit, are symmetric to *rotate left* and *link left*, respectively.

```

procedure link left;
  t, l, right(l) := right(t), t, t
end link left;

```

```

procedure rotate left;
  t, right(t), left(right(t)) := right(t), left(right(t)), t
end rotate left;

```

```

procedure assemble;
  left(r), right(l) := right(t), left(t);
  left(t), right(t) := right(null), left(null)
end assemble;

```

We can simplify top-down splaying as we did bottom-up splaying, by eliminating the second linking operation in the zig-zag case. The resulting program is as follows:

```

procedure simple top-down splay(i, t);
  if l, r: l = r = null  $\rightarrow$ 
    left(null) := right(null) := null;
  do i < item(t)  $\rightarrow$ 
    if i < item(left(t))  $\rightarrow$  rotate right
      | i  $\geq$  item(left(t))  $\rightarrow$  skip
    fi;
    link right
    | i > item(t)  $\rightarrow$ 
      if i > item(right(t))  $\rightarrow$  rotate left
        | i  $\leq$  item(right(t))  $\rightarrow$  skip
      fi;
      link left
    od{i = item(t)};
  assemble
fi
end simple top-down splay;

```

Lemma 1 holds for both top-down splaying methods with an unchanged constant factor of three. It is easy to modify the program for either method to carry out the various tree operations discussed in Section 3. For these operations a top-down method is more efficient than a bottom-up method, but the choice between the top-down methods is unclear; the simple method is easier to program but probably somewhat less efficient. Experiments are needed to discover the most practical splaying method.

5. *Semi-adjusting Search Trees*

A possible drawback of splaying is the large amount of restructuring it does. In this section we study three techniques that reduce the amount of restructuring but preserve at least some of the properties of splay trees. As in Section 4, we shall restrict our attention to successful accesses.

Our first idea is to modify the restructuring rule so that it rotates only some of the edges along an access path, thus moving the accessed node only partway toward the root. Semisplaying, our restructuring heuristic based on this idea, differs from ordinary bottom-up splaying only in the zig-zig case: after rotating the edge joining the parent $p(x)$ with the grandparent $g(x)$ of the current node x , we do not rotate the edge joining x with $p(x)$, but instead continue the splaying from $p(x)$ instead of x . (See Figure 13.)

The effect of a semisplaying operation is to reduce the depth of every node on the access path to at most about half of its previous value. If we measure the cost of a semisplaying operation by the depth of the accessed node, then Lemma 1 holds for semisplaying with a reduced constant factor of two in place of three. Furthermore, only one rotation is performed in the zig-zag case, but two steps are taken up the tree.

Like splaying, semisplaying has many variants. We describe only one, a top-down version, related to a method suggested by C. Stephenson (private communication). Think of performing a semisplaying operation as described above, bottom-up, except that if the access path contains an odd number of edges, perform the zig case at the bottom of the path instead of at the top. We can simulate this variant of semisplaying top-down, as follows. As in top-down splaying, we maintain a left tree, a middle tree, and a right tree. In addition we maintain a *top tree* and a node *top* in the top tree having a vacant child. The relationship among the trees is that all items in the left tree are less than the accessed item i and also less than those in the middle tree; all items in the right tree are greater than i and also greater than those in the middle tree; and all items in the left, middle, and right trees fall between the item in *top* and the item in its predecessor in the top tree if the vacant child of *top* is its left, or between the item in *top* and the item in its successor in the top tree if the vacant child of *top* is its right. Initially the left, right, and top trees are empty and the middle tree is the entire original tree.

Let i be the item to be accessed. Each splaying step requires looking down two steps in the middle tree from the root and restructuring the four trees according to whether these steps are to the left or to the right. If i is in the root of the middle tree, we combine the left, middle, and right trees as in the completion of top-down splaying (see Figure 12) and then make the root of the combined tree (which contains i) a child of *top*, filling its vacancy. This completes the splaying. On the other hand, if i is not in the root of the middle tree, we carry out a zig, zig-zig, or zig-zag step as appropriate.

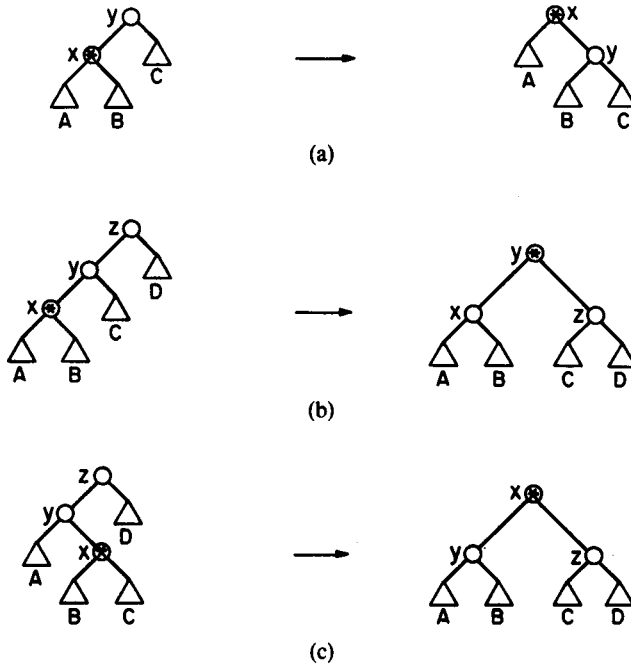


FIG. 13. A semisplaying step. Symmetric variants of cases are omitted. In each case, the starred node is the current node of the splaying operation. (a) Zig. (b) Zig-zig. (c) Zig-zag.

The zig and zig-zag cases are exactly as in top-down splaying; they do not affect the top tree. The zig-zig case is as illustrated in Figure 14. Suppose that the access path to i contains the root x of the middle tree, its left child y , and the left child of y , say z . We perform a right rotation on the edge joining x and y . Then we assemble all four trees as in the terminating case, making node y (now the root of the middle tree) a child of top and making the left and right trees the left and right subtrees of y . Finally, we break the link between z and its new parent, making the subtree rooted at z the new middle tree, the rest of the tree the new top tree, and the old parent of z the new top . The left and right trees are reinitialized to be empty.

A little thought will verify that top-down semisplaying indeed transforms the tree in the same way as bottom-up semisplaying with the zig step, if necessary, done at the bottom of the access path. Lemma 1 holds for top-down semisplaying with a constant factor of two.

Whether any version of semisplaying is an improvement over splaying depends on the access sequence. Semisplaying may be better when the access pattern is stable, but splaying adapts much faster to changes in usage. All the tree operations discussed in Section 3 can be implemented using semisplaying, but this requires using the alternative implementations of insertion and deletion. Also, the *join* and *split* algorithms become more complicated. The practical efficiency of the various splaying and semisplaying methods remains to be determined.

Another way to reduce the amount of restructuring in splay trees is to splay only sometimes. We propose two forms of conditional splaying. For simplicity we restrict our attention to a sequence of successful accesses on a tree containing a fixed set of items. One possibility is to splay only when the access path is abnormally

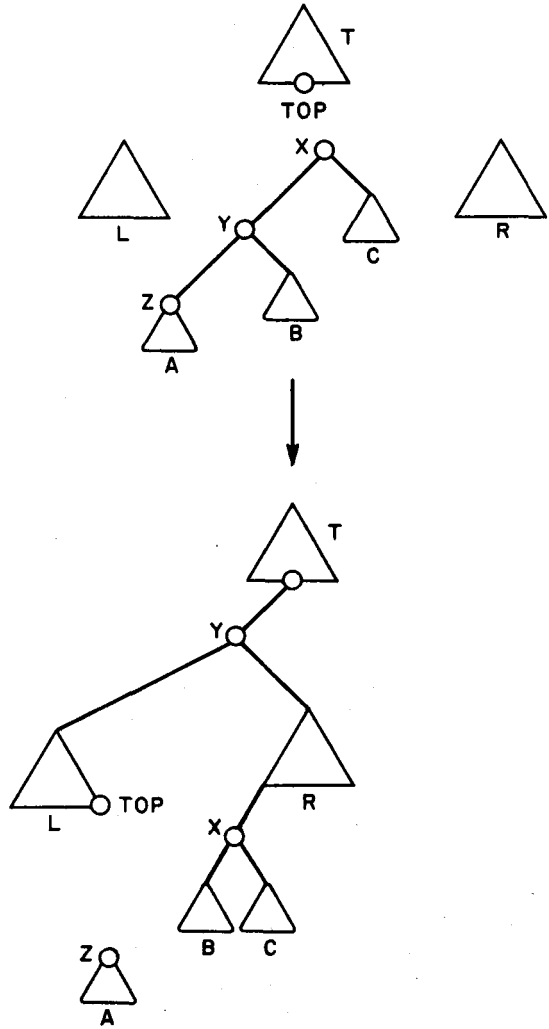


FIG. 14. Zig-zig case of a top-down semisplaying step. The symmetric variant is omitted. The initial left, right, and top trees are L , R , and T , respectively. Nodes x , y , and z are on the access path. Tree A is the new middle tree, the new left and right trees are empty, and the rest of the nodes form the new top tree.

long. To define what we mean by “abnormally long,” consider any variant of splaying (or semisplaying) such that Lemma 1 holds with a constant factor of c . Suppose every item i has a fixed weight $w(i)$, and let c' be any constant greater than c . We call an access path for item i long if the depth of the node containing i is at least $c' \log(W/w(i)) + c'/c$, and short otherwise. (Recall that W is the total weight of all the items.)

THEOREM 7 (LONG SPLAY THEOREM). *If we splay only when the access path is long and do no restructuring when the access path is short, then the total splaying time is $O(\sum_{i=1}^m \log(W/w(i)))$, that is, proportional to the amortized time to access each item once. The constant factor in this estimate is proportional to $c'/(c' - c)$. Thus the total restructuring time is independent of m , the number of accesses.*

PROOF. Consider a splaying operation on a node of depth $d \geq c' \log(W/w(i)) + c'/c$. Let ϕ and ϕ' , respectively, be the potential before and after the splaying. Since $c'/c > 1$, the actual time of the splaying is d , and we have by Lemma 1 that $d + \phi' - \phi \leq c \log(W/w(i)) + 1$. Thus the splaying reduces the potential of the tree by at least $d - c \log(W/w(i)) - 1$. This means that if we amortize the potential

