

AMORTIZED ANALYSIS OF ALGORITHMS FOR  
SET UNION WITH BACKTRACKING

Jeffery Westbrook  
Robert E. Tarjan

CS-TR-103-87

May 1987

# Amortized Analysis of Algorithms for Set Union with Backtracking\*

*Jeffery Westbrook*

Computer Science Department  
Princeton University  
Princeton, New Jersey 08544

*Robert E. Tarjan*

Computer Science Department  
Princeton, New Jersey 08544  
and  
AT&T Bell Laboratories  
Murray Hill, New Jersey 07974

May, 1987

Revised, March, 1988

## ABSTRACT

Mannila and Ukkonen have studied a variant of the classical disjoint set union (equivalence) problem in which an extra operation, called *deunion*, can undo the most recently performed union operation not yet undone. They proposed a way to modify standard set union algorithms to handle deunion operations. We analyze several algorithms based on their approach. The most efficient such algorithms have an amortized running time of  $O(\log n / \log \log n)$  per operation, where  $n$  is the total number of elements in all the sets. These algorithms use  $O(n \log n)$  space, but the space usage can be reduced to  $O(n)$  by a simple change. We prove that any separable pointer-based algorithm for the problem requires  $\Omega(\log n / \log \log n)$  time per operation, thus showing that our upper bound an amortized time is tight.

---

\* Research partially supported by National Science Foundation Grant DCR-8605962 and Office of Naval Research Contract N00014-87-K-0467.

# Amortized Analysis of Algorithms for Set Union with Backtracking\*

*Jeffery Westbrook*

Computer Science Department  
Princeton University  
Princeton, New Jersey 08544

*Robert E. Tarjan*

Computer Science Department  
Princeton, New Jersey 08544  
and  
AT&T Bell Laboratories  
Murray Hill, New Jersey 07974

May, 1987

Revised, March, 1988

## 1. Introduction.

The classical disjoint set union problem is that of maintaining a collection of disjoint sets whose union is  $U = \{1, 2, \dots, n\}$  subject to a sequence of  $m$  intermixed operations of the following two kinds:

*find* ( $x$ ): Return the name of the set currently containing element  $x$ .

*union* ( $A, B$ ): Combine the sets named  $A$  and  $B$  into a new set, named  $A$ .

---

\* Research partially supported by National Science Foundation Grant DCR-8605962 and Office of Naval Research Contract N00014-87-K-0467.

The initial collection consists of  $n$  singleton sets,  $\{1\}, \{2\}, \dots, \{n\}$ . The name of initial set  $\{i\}$  is  $i$ . For simplicity in stating bounds we assume  $m = \Omega(n)$ . This assumption does not significantly affect any of the results, and it holds in most applications.

Several fast algorithms for this problem are known [9,12]. They all combine a rooted tree set representation with some form of path compaction. The fastest such algorithms run in  $O(\alpha(m,n))$  amortized time\* per operation, where  $\alpha$  is a functional inverse of Ackermann's function [9,12]. No better bound is possible for any pointer-based algorithm that uses a separable set representation [10]. For the special case of the problem in which the subsequence of union operations is known in advance, the use of address arithmetic techniques leads to an algorithm with an amortized time bound of  $O(1)$  per operation [2].

Mannila and Ukkonen [6] studied a generalization of the set union problem called *set union with backtracking*, in which the following third kind of operation is allowed:

*deunion*: Undo the most recently performed union operation that has not yet been undone.

This problem arises in Prolog interpreter memory management [5]. Mannila and Ukkonen showed how to extend path compaction techniques to handle backtracking. They posed the question of determining the inherent complexity of the problem, and they claimed an  $O(\log \log n)$  amortized time bound per operation for one algorithm based on their approach. Unfortunately, their upper bound argument is faulty.

In this paper we derive upper and lower bounds on the amortized efficiency of algorithms for set union with backtracking. We show that several algorithms based on the approach of Mannila and Ukkonen run in  $O(\log n / \log \log n)$  amortized time per operation. These algorithms use  $O(n \log n)$  space, but the space can be reduced to  $O(n)$  by a simple change. We also show that any pointer-based algorithm that uses a separable set representation requires  $\Omega(\log n / \log \log n)$  amortized time per operation. All the algorithms we analyze are subject to this lower bound. Improving the upper

---

\* The *amortized time* is the time of an operation averaged over a worst-case sequence of operations. See the second author's survey paper [11].

bound of  $O(\log n / \log \log n)$ , if it is possible, will require the use of either a nonseparable pointer-based data structure or of address arithmetic techniques.

The remainder of this paper consists of four sections. In Section 2 we review six algorithms for set union without backtracking and discuss how to extend them to handle backtracking. In Section 3 we derive upper bounds for the amortized running times of these algorithms. In Section 4 we derive a lower bound on amortized time for all separable pointer-based algorithms for the problem. Section 5 contains concluding remarks and open problems.

## 2. Algorithms for Set Union with Backtracking

The known efficient algorithms for set union without backtracking [9,12] use a collection of disjoint rooted trees to represent the sets. The elements in each set are the nodes of a tree, whose root contains the set name. Each element contains a pointer to its parent. Associated with each set name is a pointer to the root of the tree representing the set. Each initial (singleton) set is represented by a one-node tree.

To perform union  $(A,B)$ , we make the tree root containing  $B$  point to the root containing  $A$ , or alternatively make the root containing  $A$  point to the root containing  $B$  and swap the names  $A$  and  $B$  between their respective elements. (This not only moves the name  $A$  to the right place but also makes undoing the union easy, as we shall see below.) The choice between these two alternatives is governed by a *union rule*. To perform *find*  $(x)$ , we follow the path of pointers from element  $x$  to the root of the tree containing  $x$  and return the set name stored there. In addition, we apply a *compaction rule*, which modifies pointers along the path from  $x$  to the root so that they point to nodes farther along the path.

We shall consider the following possibilities for the union and compaction rules:

### *Union Rules:*

*Union by weight:* Store with each tree root the number of elements in its tree.

When doing a union, make the root of the smaller tree point to the root of the larger, breaking a tie arbitrarily.

*Union by rank:* Store with each tree root a nonnegative integer called its rank.

The rank of each initial tree root is zero. When doing a union, make the root of smaller rank point to the root of larger rank. In the case of a tie, make either root point to the other, and increase the rank of the root of the new tree by one.

*Compaction Rules* (see Figure 1):

*Compression:* After a find, make every element along the find path point to the tree root.

*Splitting:* After a find, make every element along the find path point to its grandparent, if it has one.

*Halving:* After a find, make every other element along the find path (the first, the third, etc.) point to its grandparent, if it has one.

[Figure 1]

The two choices of a union rule and three choices of a compaction rule give six possible set union algorithms. Each of these has an amortized running time of  $O(\alpha(m,n))$  per operation [12].

We shall describe two ways to extend these and similar algorithms to handle deunion operations. The first method is the one proposed by Mannila and Ukkonen; the second is a slight variant.

We call a union operation that has been done but not yet undone *live*. We denote a pointer from a node  $x$  to a node  $y$  by  $(x,y)$ . Suppose that we perform finds without doing any compaction. Then performing deunions is easy: to undo a set union we merely make null the pointer added to the data structure by the union. To facilitate this, we maintain a *union stack*, which contains the tree roots made nonroots by live unions. To perform a deunion, we pop the top element on the union stack and make the corresponding parent pointer null.

This method works with either of the two union rules. Some bookkeeping is needed to maintain set names and sizes or ranks. Each entry on the union stack must contain not only an element but also a bit that indicates whether the corresponding union operation swapped set names. If union by rank is used, each such entry must

contain a second bit that indicates whether the union operation incremented the rank of the new tree root. The time to maintain set names and sizes or ranks is  $O(1)$  per union or deunion; thus each union or deunion takes  $O(1)$  time, worst-case. Either union rule guarantees a maximum tree depth of  $O(\log n)$  [12]; thus the worst-case time per find is  $O(\log n)$ . The space needed by the data structure is  $O(n)$ .

Mannila and Ukkonen's goal was to reduce the time per find, possibly at the cost of increasing the time per union or deunion and increasing the space. They developed the following method for allowing compaction in the presence of deunions. Let us call the forest maintained by the noncompacting algorithm described above the *reference forest*. In the compacting method, each element  $x$  has an associated *pointer stack*  $P(x)$ , which contains the outgoing pointers that have been created during the course of the algorithm but have not yet been destroyed. The bottommost pointer on this stack is one created by a union. Such a pointer is called a *union pointer*. The other pointers on the stack are ones created by compaction. They are called *find pointers*. Each pointer  $(x,y)$  of either type is such that  $y$  is a proper ancestor of  $x$  in the reference forest.

Each pointer has an *associated union operation*, which is the one whose undoing would invalidate the pointer. To be more precise, for a pointer  $(x,y)$  the associated union operation is the one that created the pointer  $(z,y)$  such that  $z$  is a child of  $y$  and an ancestor of  $x$  in the reference forest. As a special case of this definition, if  $(x,y)$  is a union pointer then  $z = x$  and the associated union operation is the one that created  $(x,y)$ . A pointer is *live* if its associated union is live.

Unions and deunions are performed as in the noncompacting method. Compactions are performed as in the set union algorithm without backtracking, except that each new pointer  $(x,y)$  is pushed onto  $P(x)$  instead of replacing the old pointer leaving  $x$ . When following a find path from an element  $x$ , the algorithm pops dead pointers from the top of  $P(x)$  until  $P(x)$  is empty or a live pointer is on top. In the former case,  $x$  is the root of its tree; in the latter case, the live pointer is followed.

This algorithm requires a way to determine whether a pointer is live or dead. For this purpose the algorithm numbers the union operations consecutively from one as they are performed. Each entry on the union stack contains the number of the corresponding union. Each pointer on a pointer stack contains the number of the

associated union and a pointer to the position on the union stack where the entry for this union was made. This information can be computed in  $O(1)$  time for any pointer  $(x,y)$  when it is created. If  $(x,y)$  is a union pointer, the information is computed as part of the union. If  $(x,y)$  is a find pointer, then the last pointer on the find path from  $x$  to  $y$  when  $(x,y)$  was created has the same associated union as  $(x,y)$  and has stored with it the needed information. To test whether a pointer is live or dead, it is merely necessary to access the union stack entry whose position is recorded with the pointer and test first, if the entry is still on the stack, and second, whether its union number is the same as that stored with the pointer. If so, the pointer is live; if not, dead.

The implementation of deunion must be changed slightly, to preserve the invariant that in every pointer stack all the dead pointers are on top. To perform a deunion, the algorithm pops the top entry on the union stack. Let  $x$  be the element in this entry. The algorithm pops  $P(x)$  until it contains only one pointer, which is the union pointer created by the union that is to be undone. The algorithm restores the set names and sizes or ranks as necessary, and pops the last pointer from  $P(x)$ . Because of the compaction, the state of the data structure after a deunion will not in general be the same as its state before the corresponding union.

We call this method the *lazy method* since it destroys dead pointers in a lazy fashion. Either of the union rules and any of the compaction rules can be used with the method. The total running time is proportional to  $m$  plus the total number of pointers created. (With any of the compaction rules, a compaction of a find path containing  $k \geq 2$  pointers results in the creation of  $\Omega(k)$  pointers,  $k - 1$  in the case of compression or splitting and  $\lfloor k/2 \rfloor$  in the case of halving.)

An alternative to the lazy method is the *eager method*, which pops pointers from pointer stacks as soon as they become dead. To make this popping possible, each union stack entry must contain a list of the pointers whose associated union is the one corresponding to the entry. When a union stack entry is popped, all the pointers on its list are popped from their respective pointer stacks as well. Each such pointer will be on top of its stack when it is to be popped. To represent such a pointer, say  $(x,y)$ , in a union stack entry, it suffices to store  $x$ . With this method, numbering the union operations is unnecessary, as is popping pointer stacks during finds.

The time required by the eager method for any sequence of operations is only a



constant factor greater than that required by the lazy method, since both methods create the same pointers but the eager method destroys them earlier. With either union rule, the eager method uses  $O(n \log n)$  space in the worst case, since the maximum tree depth is  $O(\log n)$  and all pointers on any pointer stack point to distinct elements. (From bottom to top, the pointers on  $P(x)$  point to shallower and shallower ancestors of  $x$ .)

The lazy method also has an  $O(n \log n)$  space bound, as observed by Esa Helttula (private communication, 1988). For any node  $x$ , consider the top pointer on  $P(x)$ , which is to a node, say  $y$ . Even if the pointer from  $x$  to  $y$  is currently dead, it must once have been live, and all pointers currently on  $P(x)$  point to distinct nodes on the tree path from  $x$  to  $y$  as it existed when the pointer from  $x$  to  $y$  was live. Thus there can be only  $O(\log n)$  such pointers.

The choice between the lazy and eager methods is not clear-cut. The lazy method requires numbering the unions, and these numbers can grow arbitrarily large, although reuse of such numbers reduces the number of distinct numbers required to  $O(n \log n)$ . As we shall see at the end of Section 3, a small change in the compaction rules reduces the space needed by either method to  $O(n)$ .

### 3. Upper Bounds on Amortized Time

The analysis to follow applies to both the lazy method and the eager method. Ignoring the choice between lazy and eager pointer deletion, there are six versions of the algorithm, depending on the choice of a union rule and a compaction rule.

As a first step on the analysis, we note that compression with either union rule is no better in the amortized sense than doing no compaction at all, i.e. the amortized time per operation is  $\Omega(\log n)$ . The following class of examples shows this. For any  $k$ , form a tree of  $2^k$  elements by doing unions on pairs of elements, then on pairs of pairs, and so on. This produces a tree called a *binomial tree*  $B_k$ , whose depth is  $k$ . (See Figure 2.) Repeat the following three operations any number of times: do a find on the deepest element in  $B_k$ , undo the most recent union, and redo the union. Each find creates  $k - 1$  pointers, which are all immediately made dead by the subsequent deunion. Thus the amortized time per operation is  $\Omega(k) = \Omega(\log n)$ .

[Figure 2]

Both splitting and halving perform better; each has an  $O(\log n / \log \log n)$  amortized bound per operation, in combination with either union rule. To prove this, we need a definition. For an element  $x$ , let  $size(x)$  be the number of descendants of  $x$  (including itself) in the reference forest. The *logarithmic size* of  $x$ ,  $lgs(x)$ , is  $\lfloor \lg size(x) \rfloor^*$ .

We need the following lemma concerning logarithmic sizes when union by weight is used.

**Lemma 1 [9].** *Suppose union by weight is used. If node  $v$  is the parent of node  $w$  in the reference forest, then  $lgs(w) < lgs(v)$ . Any node has logarithmic size between 0 and  $\lg n$  (inclusive).*

**Proof.** When a node  $v$  becomes the parent of another node  $w$ ,  $size(w) \leq 2 size(v)$  by the union by weight rule. Later unions can only increase  $size(v)$  and cannot increase  $size(w)$  (unless the union linking  $v$  and  $w$  is undone). The lemma follows.  $\square$

**Theorem 1.** *Union by weight in combination with either splitting or halving gives an algorithm for set union with backtracking running in  $O(\log n / \log \log n)$  amortized time per operation.*

**Proof.** We shall charge the pointer creations during the algorithm to unions and finds in such a way that each operation is charged for  $O(\log n / \log \log n)$  pointer creations. For an arbitrary positive constant  $c < 1$ , we call a pointer  $(x,y)$  *short* if  $lgs(y) - lgs(x) \leq c \lg \lg n$  and *long* otherwise. (The logarithmic sizes in this definition are measured at the time  $(x,y)$  is created.) We charge the creation of a pointer  $(x,y)$  as follows:

- (i) If  $y$  is a tree root, charge the operation (union or find) that created  $(x,y)$ .
- (ii) If  $y$  is not a tree root and  $(x,y)$  is long, charge the find that created  $(x,y)$ .
- (iii) If  $y$  is not a tree root and  $(x,y)$  is short, charge the union that most recently made

---

\* For any  $x$ ,  $\lg x = \log_2 x$ .

$y$  a non-root.

A find with splitting creates two new paths of pointers, and a find with halving creates one new path of pointers. Thus  $O(1)$  pointers are charged to each operation by (i). The number of long pointers along any path can be estimated as follows. For any long pointer  $(x,y)$ ,  $lgs(y) - lgs(x) > c \lg \lg n$ . Logarithmic sizes strictly increase along any path and are between 0 and  $\lg n$  by Lemma 1. Thus if there are  $k$  long pointers on a path,  $\lg n \geq k c \lg \lg n$ , which implies  $k \leq \lg n / (c \lg \lg n)$ . Thus a find with either splitting or halving can create only  $O(\log n / \log \log n)$  long pointers, which means that  $O(\log n / \log \log n)$  pointers are charged to each find by (ii).

It remains for us to bound the number of pointers charged by (iii). Consider a union operation that makes an element  $x$  a child of another element  $y$ . Let  $I$  be the time interval during which pointers are charged by (iii) to this union. During  $I$ , the sizes, and hence the logarithmic sizes, of all descendants of  $x$  remain constant. Interval  $I$  ends with the undoing of the union.

For each descendant  $w$  of  $x$ , at most one pointer  $(w,x)$  can be charged by (iii) to the union, since the creation of another such pointer charged by (iii) cannot occur at least until  $x$  again becomes a root and then becomes a nonroot, which can only happen after the end of  $I$ . Thus the number of pointers charged by (iii) to the union is at most one per descendant  $w$  of  $x$  such that  $lgs(x) - lgs(w) \leq c \lg \lg n$ .

Since logarithmic sizes strictly increase along tree paths, any two elements  $u$  and  $v$  with  $lgs(u) = lgs(v)$  must be unrelated, i.e. their sets of descendants are disjoint. This means that the number of descendants  $w$  of  $x$  with  $lgs(w) = i$  is at most  $size(x)/2^i \leq 2^{lgs(x)+1-i}$ , and the number of descendants  $w$  of  $x$  with  $lgs(x) - lgs(w) \leq c \lg \lg n$  is at most

$$\sum_{i=lgs(x)-\lfloor c \lg \lg n \rfloor}^{lgs(x)} 2^{lgs(x)+1-i} \leq 2^{\lfloor c \lg \lg n \rfloor + 2} = O((\log n)^c) = O(\log n / \log \log n),$$

since  $c < 1$ . Thus there are  $O(\log n / \log \log n)$  pointers charged to the union by (iii).  
□

The same result holds if union by rank is used instead of union by weight, but in this case the proof becomes a little more complicated because logarithmic sizes need

not strictly increase along tree paths. We deal with this by slightly changing the definition of short and long pointers. We need the following lemma.

**Lemma 2 [12].** *Suppose union by rank is used. If node  $v$  is the parent of node  $w$  in the reference forest, then  $0 \leq lgs(w) \leq lgs(v) \leq \lg n$  and  $0 \leq rank(w) < rank(v) \leq \lg n$ .*

**Proof.** The first group of inequalities is immediate. The definition of union by rank implies  $rank(w) < rank(v)$ . A proof by induction on the rank of  $v$  shows that  $size(v) \geq 2^{rank(v)}$ , which implies that  $rank(v) \leq \lg n$ .  $\square$

**Theorem 2.** *Union by rank in combination with either splitting or halving gives an algorithm for set union with backtracking running in  $O(\log n / \log \log n)$  amortized time per operation.*

**Proof.** We define a pointer  $(x,y)$  to be *short* if  $\max\{lgs(y) - lgs(x), rank(y) - rank(x)\} \leq c \lg \lg n$  and *long* otherwise, where  $c < 1$  is a positive constant. We charge the creation of pointers to unions and finds exactly as in the proof of Theorem 1 (rules (i), (ii), and (iii)). The number of pointers charged by rule (i) is  $O(1)$  per union or find, exactly as in the proof of Theorem 1. A long pointer  $(x,y)$  satisfies at least one of the inequalities  $lgs(y) - lgs(x) > c \lg \lg n$  and  $rank(y) - rank(x) > c \lg \lg n$ . Along any tree path only  $O(\log n / \log \log n)$  long pointers can satisfy the former inequality and only  $O(\log n / \log \log n)$  long pointers can satisfy the latter, by Lemma 2. It follows that only  $O(\log n / \log \log n)$  pointers can be charged per find by rule (ii).

To count short pointers, we make one more definition. For a non-root element  $x$ , let  $p(x)$  be the parent of  $x$  in the reference forest. A non-root  $x$  is *good* if  $lgs(x) < lgs(p(x))$  and *bad* otherwise, i.e. if  $lgs(x) = lgs(p(x))$ . The definition of *lgs* implies that any element can have at most one bad child. The bad elements thus form paths called *bad paths* of length  $O(\log n)$ ; all elements on a bad path have the same logarithmic size. We call the element of largest rank on a bad path the *head* of the path. The head of a bad path is a bad element whose parent is either a good element or a tree root.

Consider a union operation that makes an element  $x$  a child of an element  $y$ . We

count short pointers charged to this union as follows:

- (1) *Short pointers leading from good elements.* If  $v$  and  $w$  are good elements such that  $lgs(v) = lgs(w)$ , then  $v$  and  $w$  are unrelated in the reference forest, i.e. they have disjoint sets of descendants. The analysis that yielded the count of short pointers in the proof of Theorem 1 applies to the good elements here to yield a bound of  $O((\log n)^c) = O(\log n / \log \log n)$  short pointers leading from good elements that are charged to the union by (iii).
- (2) *Short pointers leading from bad elements.* Consider the number of bad paths from which short pointers can lead to  $x$ . The head of such a path is an element  $w$  such that  $p(w)$  is either good or a tree root, and  $lgs(x) - lgs(w) \leq c \lg \lg n$ . Heads of different bad paths have different parents. The analysis that counts short pointers in the proof of Theorem 1 yields an  $O((\log n)^c)$  bound on the number of bad paths from which short pointers can lead to  $x$ . Along such a bad path, rank strictly increases, and the definition of shortness implies that only the  $c \lg \lg n$  elements of largest rank along the path can have short pointers leading to  $x$ . The total number of short pointers leading from bad nodes that are charged to the union by (iii) is thus  $O(c \log \log n (\log n)^c) = O(\log n / \log \log n)$ .  $\square$

We conclude this section by discussing how to reduce the space bound for both the lazy method and the eager method to  $O(n)$ . This is accomplished by making the following simple changes in the compaction rules. If union by size is used, the compaction of a find path is begun at the first node along the path whose size is at least  $\lg n$ . If union by rank is used, the compaction of a find path is begun at the first node whose rank is at least  $\lg \lg n$ . With this modification, only  $O(n / \log n)$  nodes have find pointers leaving them, and the total number of pointers in the data structure at any time is  $O(n)$ . The analysis in Theorems 1 and 2 remains valid, except that there is an additional time per find of  $O(\log \log n)$  to account for the initial, noncompacted part of each find path.

#### 4. A General Lower Bound on Amortized Time

We shall prove that the bound in Theorems 1 and 2 is best possible for a large class of algorithms for set union with backtracking. Our computational model is the *pointer machine* [3,4,8,10 ] with an added assumption about the data structure called *separability*. Related results follow. Tarjan [10] derived an amortized bound in this model for the set union problem without backtracking. Blum [1] derived a worst-case-per-operation lower bound for the same problem. Mehlhorn, Näher, and Alt [7] derived an amortized lower bound for a related problem. Their result does not require separability.

The algorithms to which our lower bound applies are called *separable pointer algorithms*. Such an algorithm uses a linked data structure that can be regarded as a directed graph, with each pointer represented by an edge. The algorithm solves the set union with backtracking problem according to the following rules:

- (i) The operations are presented on-line, i.e. each operation must be completed before the next one is known.
- (ii) Each set element is a node of the data structure. There can be any number of additional nodes.
- (iii) (Separability). After any operation, the data structure can be partitioned into node-disjoint subgraphs, one corresponding to each currently existing set and containing all the elements in the set. The name of the set occurs in exactly one node in the subgraph. *No edge leads from one subgraph to another.*
- (iv) The cost of an operation *find* ( $x$ ) is the length (number of edges) of the shortest path from  $x$  to the node that holds the name of the set containing  $x$ . This length is measured at the beginning of the find, i.e. before the algorithm changes the structure as specified in (v).
- (v) During any *find*, *union*, or *deunion* operation, the algorithm can add edges to the data structure at a cost of one per edge, delete edges at a cost of zero, and move, add, or delete set names at a cost of zero. The only restriction is that separability must hold after each operation.

The eager method of Section 2 obeys rules (i)-(v). This is also true of the lazy method, if we regard pointers as disappearing from the model data structure as soon as they become dead. This does not affect the performance of the algorithm in the model, since once a pointer becomes dead it is never followed.

**Theorem 3.** *For any  $n$ , any  $m = \Omega(n)$ , and any separable pointer algorithm, there is a sequence of  $m$  find, union, and deunion operations whose cost is  $\Omega(m \log n / \log \log n)$ .*

**Proof.** We shall prove the theorem for  $n$  of the form  $2^{2^k}$  for some  $k \geq 1$  and for  $m \geq 4n$ . The result follows for all  $n$  and  $m = \Omega(n)$  by padding the expensive problem instances constructed below with extra singleton sets on which no operations take place and with extra finds.

In estimating the cost of a sequence of operations, we shall charge the cost of adding an edge to the data structure to the deletion of the edge. Since this postpones the cost, it cannot increase the total cost of a sequence.

We construct an expensive sequence as follows. The first  $n - 1$  operations are unions that build a set of size  $n$  by combining singletons in pairs, pairs in pairs, and so on. The remaining operations occur in groups, each group containing between 1 and  $2n - 2$  operations. Each group begins and ends with all the elements in one set. We obtain a group of operations by applying the appropriate one of the following two cases (if both apply, either may be selected). Let  $b = \lfloor \lg n / (2 \lg \lg n) \rfloor$ .

- (1) If some element in the (only) set is at distance at least  $b$  away from the set name, do a find on this element.
- (2) If some sequence of  $\ell$  deunions will force the deletion of  $\ell b$  edges from the data structure (to maintain separability), do these deunions. Then do the corresponding unions in the reverse order, restoring the initial set of size  $n$ .

We claim that if there is only one set, formed by repeated pairing, then Case (1) or Case (2) must apply. If this is true, we can obtain an expensive sequence of operations by generating successive groups of operations until more than  $m - 2n + 2$  operations have occurred, and then padding the sequence with enough additional finds to

make a total of  $m$  operations. The cost of such a sequence is at least  $(m-3n+3)b = \Omega(m \log n / \log \log n)$ .

It remains to prove the claim. Suppose Case (2) does not apply. We shall show that Case (1) does. Let  $f = (\lg n)^2$ . For  $0 \leq i \leq \lg n / \lg f$  we define a partition  $P_i$  of the nodes of the data structure by

$P_i = \{X \mid X \text{ is the collection of nodes in the subgraph corresponding to one of the sets that would be formed by doing } f^i - 1 \text{ deunions}\}$

Observe that  $|P_i| = f^i$ . Also  $f^{\lg n / \lg f} = n$ , so  $P_i$  is defined for  $i \leq \lg n / \lg f$ . In particular  $P_b$  is defined, since  $b = \lfloor \lg n / (2 \lg \lg n) \rfloor = \lfloor \lg n / \lg f \rfloor$ .

For  $0 \leq i \leq \lg n / \lg f$ , we define the collection  $D_i$  of *deep sets* in  $P_i$  by

$D_i = \{X \in P_i \mid \text{all elements in } X \text{ are at distance at least } i \text{ from the name of the single set}\}$ .

Let  $d_i = |D_i|$ . We shall show that  $d_b > 0$ , which implies the existence of an element at distance at least  $b$  away from the name of the single set; hence Case (1) applies.

Let  $\ell_i$  be the number of edges that lead from one set in  $P_i$  to another. We have  $\ell_i \leq bf^i$ , since otherwise performance of  $f^i - 1$  deunions would force the deletion of  $bf^i$  edges from the data structure, and Case (2) would apply.

Now we derive a recursive bound on  $d_i$ . We have  $d_1 = f - 1$ , since only one of the  $f$  sets in  $P_1$  can contain the only set name. We claim that  $d_{i+1} \geq fd_i - \ell_i$ . To verify the claim, let us consider  $D_i$ . Since  $n = 2^{2^k}$  and the union structure of the only set forms a binomial tree, each set  $X$  in  $D_i$  consists of  $f$  sets in  $P_{i+1}$ , all of whose elements are at distance at least  $i$  from the name of the only set. For an element  $x \in X$  to be at distance exactly  $i$  from the set name, some edge must lead from  $x$  to a set in  $P_i$  other than  $X$ ; otherwise  $X$  would not be in  $D_i$ . There are  $\ell_i$  such edges. Each such edge can eliminate one set in  $P_{i+1}$  from being in  $D_{i+1}$ . But this leaves  $fd_i - \ell_i$  sets in  $D_{i+1}$ , namely the  $fd_i$  sets into which the sets in  $D_i$  divide, minus at most  $\ell_i$  eliminated by edges between different sets in  $P_i$ . That is,  $d_{i+1} \geq fd_i - \ell_i$ , as claimed.

Applying the bound  $\ell_i \leq bf^i$  gives  $d_{i+1} \geq fd_i - bf^i$ . Using  $d_1 = f - 1$ , a proof by induction shows that  $d_i \geq f^{i-1}(f - (i-1)b - 1)$ .



We wish to show that  $d_b > 0$ . This is true provided that  $(f - (b-1)b - 1) > 0$ . But  $f = (\lg n)^2$  and  $b = \lfloor \lg n / (2 \lg \lg n) \rfloor$ , giving  $(f - (b-1)b - 1) = (f - b^2 + b - 1) \geq \frac{3}{4} (\lg n)^2 > 0$ , since we are assuming  $n \geq 4$ , which implies  $b^2 \leq (\lg n)^2 / 4$  and  $b \geq 1$ . Thus  $d_b > 0$ , which implies that some element is at distance at least  $b$  from the set name, i.e. Case (1) applies.  $\square$

## 5. Remarks

Our bound of  $\Theta(\log n / \log \log n)$  on the amortized time per operation in the set union problem with backtracking is the same as Blum's worst-case bound per operation in the set union problem without backtracking [1]. Perhaps this is not a coincidence. Our lower bound proof resembles his. Furthermore the data structure he uses to establish his upper bound can easily be extended to handle deunion operations; the worst-case bound per operation remains  $O(\log n / \log \log n)$  and the space needed is  $O(n)$ .

The compaction methods have the advantage over Blum's method that as the ratio of finds to unions and deunions increases the amortized time per find decreases. The precise result is that if the ratio of finds to unions and deunions in the operation sequence is  $\gamma$  and the amortized time per union and deunion is defined to be  $\Theta(1)$ , then the amortized time per find is  $\Theta(\log n / (\max\{1, \log(\gamma \log n)\}))$ . This bound is valid for any value of  $\gamma$ , and it holds for splitting or halving with either union rule, and it is the best bound possible for any separable pointer algorithm. This can be proved using straightforward extensions of the arguments in Sections 3 and 4. The space bound can be made  $O(n)$  by an extension of the idea proposed at the end of Section 3. If the deunion operations occur in bursts, the time per operation decreases further, but we have not attempted to analyze this situation.

Perhaps the most interesting open problem is whether the lower bound in Section 4 can be extended to nonseparable pointer algorithms. (In place of separability, we require that the out-degree of every node in the data structure be constant.) We conjecture that the bound in Theorem 3 holds for such algorithms. The techniques of Mehlhorn, Näher, and Alt [7] suggest an approach to this question, which might yield at least an  $\Omega(\log \log n)$  bound if not an  $\Omega(\log n / \log \log n)$  bound on the amortized

time.

## References

- [1] N. Blum, "On the single-operation worst-case time complexity of the disjoint set union problem," *SIAM J. Comput.* 15 (1986), 1021-1024.
- [2] H. N. Gabow and R. E. Tarjan, "A linear-time algorithm for a special case of disjoint set union," *J. Comput. Sys. Sci.* 30 (1985), 209-221.
- [3] D. E. Knuth, *The Art of Computer Programming, Vol. 1, Fundamental Algorithms*, Addison-Wesley, Reading, MA, 1968.
- [4] A. N. Kolmogorov, "On the notion of algorithm," *Uspehi Mat. Nauk.* 8 (1953), 175-176.
- [5] H. Mannila and E. Ukkonen, "On the complexity of unification sequences," *Third International Conference on Logic Programming*, July 14-18, 1986, *Lecture Notes in Computer Science* 225, Springer-Verlag, New York, 1986, 122-133.
- [6] H. Mannila and E. Ukkonen, "The set union problem with backtracking," *Proc. Thirteenth International Colloquium on Automata, Languages, and Programming (ICALP 86)*, Rennes, France, July 15-19, 1986, *Lecture Notes in Computer Science* 226, Springer-Verlag, New York, 1986, 236-243.
- [7] K. Mehlhorn, S. Näher, and H. Alt, "A lower bound for the complexity of the union-split-find problem," *Proc. Fourteenth International Colloquium on Automata Languages, and Programming (ICALP 87)*, Karlsruhe, West Germany, July 13-17, 1987, *Lecture Notes in Computer Science*, Springer-Verlag, to appear.
- [8] A. Schönhage, "Storage modification machines," *SIAM J. Comput.* 9 (1980), 490-508.
- [9] R. E. Tarjan, "Efficiency of a good but not linear set union algorithm," *J. Assoc. Comput. Mach.* 22 (1975), 215-225.
- [10] R. E. Tarjan, "A class of algorithms which require nonlinear time to maintain disjoint sets," *J. Comput. Sys. Sci.* 18 (1979), 110-127.
- [11] R. E. Tarjan, "Amortized computational complexity," *SIAM J. Alg. Disc. Meth.* 6 (1985), 306-318.

- [12] R. E. Tarjan and J. van Leeuwen, "Worst-case analysis of set union algorithms,"  
*J. Assoc. Comput. Mach.* 31 (1984), 245-281.

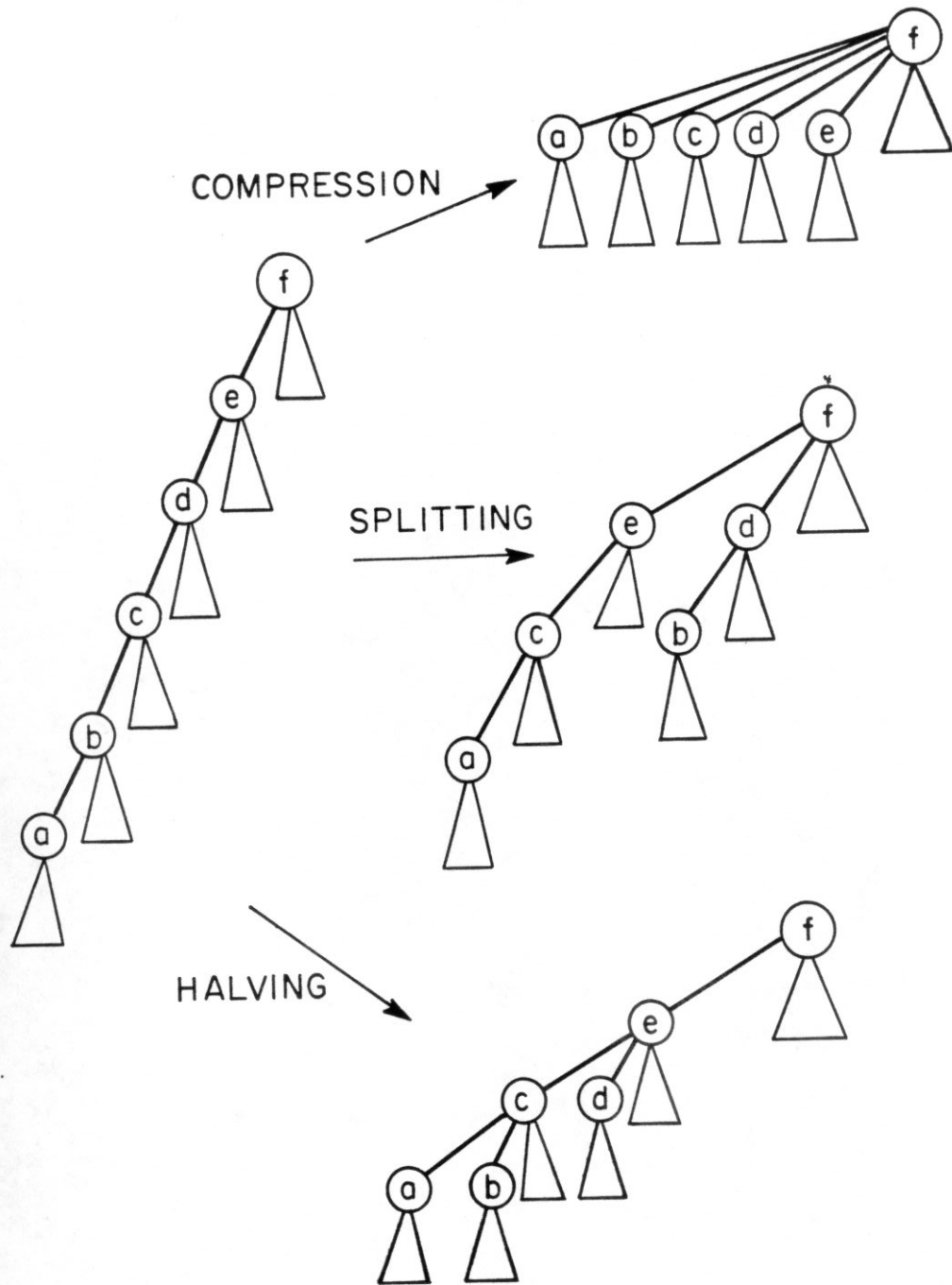


Figure 1. Path compression, path splitting, and path halving. The element found is "a".

$$B_0 = \circ$$

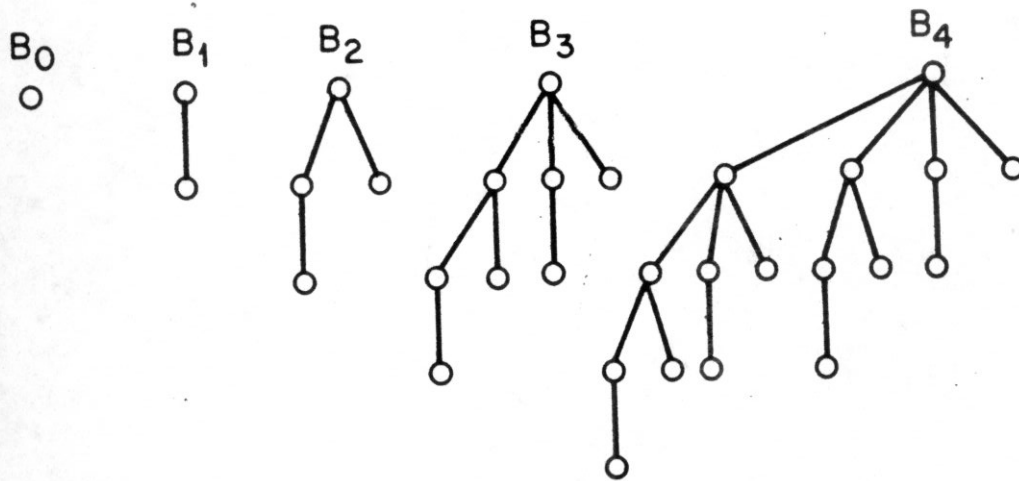
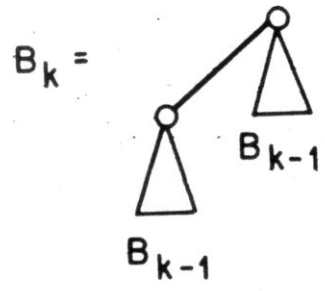


Figure 2. Binomial trees.