Multilevel Data Structures: Models and Performance

ABHA MOITRA, S. SITHARAMA IYENGAR, FAROKH B. BASTANI, MEMBER, IEEE, AND I. LING YEN

Abstract—We advocate a stepwise method of deriving high performance implementation of a set of operations. This method is based on the ability to organize the data into a multilevel data structure so as to provide an efficient implementation of all the operations. Typically, for such data organization the performance may deteriorate over a period of time and that can be corrected by reorganizing the data. This data reorganization is done by the introduction of maintenance processes.

For a particular example we consider the multilevel data organization and the different models of maintenance processes possible. The various models of maintenance processes provide varying amounts of concurrency by varying the degree of atomicity in different operations. Performance behavior for the different models are derived and we sketch a correctness proof for the developed implementation. Simulation studies of the performance for this example confirm that the performance improves as we move from coarse grained concurrency to finer grained concurrency.

Index Terms-Maintenance processes, multilevel data structures, performance, program correctness, weak/strong invariants.

I. INTRODUCTION

S YSTEM performance can be improved by scheduling house keeping activities to occur during periods when the processor is waiting for messages or is transmitting messages. In other words performance of a distributed system can be improved by scheduling the computations required to maintain the data structure during the idle periods. This approach is particularly suitable for distributed systems since the data reorganization can be done while a response is being sent to the invoker or while the local processor is waiting for a new request.

For example, it is efficient to search for a key in a perfectly balanced binary search tree. However, rebalancing the tree after every insertion/deletion operation is inefficient. A maintenance process in a distributed system can rebalance the tree while a response is being sent to the invoker or when the processor is waiting for a new request. This method is proposed by Lampson [8] and Manber [9]. Manber [9] gives the algorithms and proofs for maintenance processes for a special form of binary search trees data structure called "external trees" and in Manber

Manuscript received September 30, 1985. The work of F. B. Bastani was supported in part by the National Science Foundation under Grant MCS-83-01745.

A. Moitra is with the Department of Computer Science, Cornell University, Ithaca, NY 14853.

S. S. Iyengar is with the Department of Computer Science, Louisiana State University, Baton Rouge, LA 70803.

F. B. Bastani is with the Department of Computer Science, University of Houston, Houston, TX 77004.

I. L. Yen is with the Department of Computer and Information Science, University of Massachusetts, Amherst, MA 01003.

IEEE Log Number 8820972.

and Ladner [10] a binary search structure is considered. In Ellis [4], a distributed version of an extendible hash file useful in distributed data bases is developed by "distributing a data structure". Bastani *et al.* [1], [2] have investigated the improvement of the performance of data abstractions in a distributed system through the use of low priority maintenance process. They consider the case where the interface operations are implemented by foreground processes while the maintenance tasks are performed by background processes. They give algorithms for these processes in a stepwise manner, proceeding from coarse grained concurrency to fine grained concurrency.

In this paper we discuss one approach for improving the average response time for a set of operations. This approach is based on the usage of multilevel data structures along with maintenance processes [1]–[3], [8], [9]. We illustrate this approach by considering one example in detail. For this example we present different models which capture the types of interactions between the processes providing the original operations and the maintenance processes. Even though we take a specific example for illustration purposes, the different models of interactions are general and can be modified to deal with other problems. Finally, we also present some simulation studies.

II. DATA STRUCTURE PERFORMANCE TRADEOFFS

The problem we consider is that of implementing a SearchTable which should provide the following operations:

- ADD(k,e): Add element e with key k to the Search-Table.
- DELETE(k): Delete the element associated with key k.
- LOOKUP(k): Return the element associated with key k.

Standard implementations cannot simultaneously optimize all the above three operations. For example,

1) a linked list is poor for LOOKUP, but good for ADD and DELETE

2) an unsorted list is poor for LOOKUP and DELETE, but good for ADD

3) a sorted list is good for LOOKUP, but poor for ADD and DELETE

4) a binary search tree has a much better performance for ADD and DELETE than a sorted array, although it has a slightly worse LOOKUP performance

5) a balanced tree improves on the performance time

0098-5589/88/0600-0858\$01.00 © 1988 IEEE

for LOOKUP at the expense of that of ADD and DE-LETE

6) hash tables with open addressing are not suitable for DELETE while those with chaining have degraded performance as the chained lists get longer.

III. MULTILEVEL DATA STRUCTURES

Typically, the data on which the operations are performed is organized in some particular data structure. However, in the multilevel data structure approach the data is organized in a number of different data structures so as to provide efficient implementation for all the operations. As a sequence of operations are performed, the organization of the data is altered. Now since the different data structures have different performance characteristics, the performance of the operations may deteriorate as a result of the new organization of the data. Consequently, some maintenance processes [9] have to be invoked to reorganize the data so as to improve the performance of the operations.

Since, an operation may be performed under different data organization, the implementation of the operation should be consistent with all the different data organizations. This is achieved by defining two types of invariants for the data organization: a strong and weak invariant, such that the following hold.

1) If the strong invariant is true then any operation can be implemented efficiently. The processing of an operation may falsify the strong invariant.

2) If the weak invariant holds then execution of an operation will maintain the weak invariant.

The maintenance process when invoked attempts to establish the strong invariant given that the weak invariant holds. Whether the maintenance process actually succeeds in this attempt depends on the model considered; however, for any model if no new requests arrive from the client for a certain period of time then the maintenance process will establish the strong invariant. This point will become clearer when we discuss the different models.

For the data type SearchTable and its associated operations, a multilevel data structure consisting of a binary search tree and a sorted array can be used as given in Fig. 1. Elements to be added are added to the binary search tree and for efficient implementation of DELETE, each element has a flag (with values "dead" and "alive") associated with it. The strong invariant is that the binary search tree is empty and that all the elements in the array have flag value "alive". One or more maintenance processes (which we will also call as background processes) transfer elements from the binary search tree to the sorted array and also remove dead elements.

The reason for calling such a data organization as a multilevel data structure is that it is similar to the memory hierarchy used in computer systems. The binary search tree in the above example can be viewed as a cache storage while the sorted array corresponds to long term storage. One possible application of such a technique is in the



Fig. 1. Multilevel data structure for the SearchTable module.

implementation of a password file for an authorization server in a local area network.

One of the earliest systematic methods of developing interfering programs is the stepwise development method of Jones [6] which uses rely/guarantee conditions for each of the processes. A general method for showing noninterference is to identify an invariant and show that execution of each statement preserves the invariant [11]. A basic approach for implementing such programs is to start from programs with coarse grained concurrency and refine them to ones with fine grained concurrency [3].

In [1] the "coarse grained concurrency to fine grained concurrency" approach is combined with the stepwise approach of [6]. This allows modular proofs of the different procedures to be undertaken and hence the addition of new operations only requires proofs about the added code.

The stepwise development from coarse grained concurrency to fine grained concurrency proceeds in the following sequence.

1) In Case A, only the foreground process is present, the data structure used is a binary search tree.

2) In Case B, background processes are added but all the processes are nonpreemptible. The multilevel data structure of Fig. 1 is used.

3) In Case C, a background process is preemptible by a foreground process. The multilevel data structure of Fig. 1 is used.

IV. MODELS FOR MULTILEVEL DATA STRUCTURES

We now develop formal models (Kleinrock [7]) for the three cases and derive their performance behavior (see also [5]). We assume that the client (i.e., the program using the abstract data type) repeatedly executes the following cycles:

Compute for a period having the exponential distribution with parameter λ ; then invoke an operation of the abstract data type.

If the average response time is R, then the productive work PW performed by the client is

$$PW = \frac{1/\lambda}{1/\lambda + R}$$

This is the proportion of time that the client is doing some other (presumably more useful) chore than waiting for an

L

operation of the abstract data type to complete. Clearly, if $1/\lambda$ is large then this abstract data type is rarely used by the client, so that optimization is not critical. However, if $1/\lambda$ is small then operations on the abstract data type could be a bottleneck.

A. Case A-No Background Processes

The state transition diagram is shown in Fig. 2. In this case the client is blocked until the current update is over.

In state s_0 , the client is doing some other work. In state s_{11} , a LOOKUP operation is being performed. In state s_{12} , an update operation (to complete the ADD or DELETE operation) is in progress. The probability that a client request needs to update the data structure is p'. The search operation is completed with rate μ_0 , while the update operation has rate μ_1 . If P_i denotes the probability of being in state s_i then the following equations are obtained for the states.

for
$$s_{12}$$
: $p'\mu_0 P_{11} = \mu_1 P_{12}$
for s_{11} : $\mu_0 P_{11} = \lambda P_0$
 $P_0 + P_{11} + P_{12} = 1.$

Solving these we obtain

$$P_0 = \frac{1/\lambda}{1/\lambda + 1/\mu_0 + p'/\mu_1}$$

Now $PW = P_0$ and hence we obtain

$$PW = \frac{1/\lambda}{1/\lambda + 1/\mu_0 + p'/\mu_1}$$
$$R = \frac{1}{\mu_0} + \frac{p'}{\mu_1}.$$

Remarks:

1) As λ increases, *PW* decreases since this module then becomes a bottleneck.

2) If the search algorithm is improved, then μ_0 increases and hence *PW* increases.

3) If there are fewer ADD/DELETE requests then p' decreases and hence *PW* increases.

4) If the update algorithm is improved then μ_1 increases and hence *PW* increases.

B. Case B-Nonpreemptible Background Processes

The state transition diagram is shown in Fig. 3. Here the client is blocked until the previous update is over.

In this implementation, the foreground process returns control to the client process as soon as the LOOKUP portion of the operation is completed. The background processes do the actual update in state s_{12} . Since there is no preemption, these processes complete the update even if the client issues a new operation. This is indicated by the transition from state s_{12} to s_2 . For this case we obtain the following equations:



Fig. 2. State transition for Case A.



Fig. 3. State transition for Case B.

for
$$s_2$$
: $\lambda P_{12} = \mu_1 P_2$
for s_{12} : $(\lambda + \mu_1) P_{12} = p' \mu_0 P_{11}$
for s_{11} : $\mu_0 P_{11} = \lambda P_0 + \mu_1 P_2$
 $P_0 + P_{11} + P_{12} + P_2 = 1$.

1

Solving these, we obtain

$$P_{0} = \frac{1}{1 + [\lambda(\lambda + \mu_{1})/(p\lambda + \mu_{1})][1/\mu_{0} + p'/\mu_{1}]}$$

$$PW = P_{0} + P_{12}$$

$$= \frac{1/\lambda}{1/\lambda + 1/\mu_{0} + [\lambda/(\lambda + \mu_{1})]p'/\mu_{1}}$$

$$\geq PW \text{ for Case A} \quad \text{since } \frac{\lambda}{\lambda + \mu_{1}} \leq 1$$

$$R = \frac{1}{\mu_{0}} + \frac{\lambda p'}{(\lambda + \mu_{1})\mu_{1}}$$

$$\leq \frac{1}{\mu_{0}} + \frac{p'}{\mu_{1}} \quad \text{since } \frac{\lambda}{\lambda + \mu_{1}} \leq 1$$

= response time for Case A.

Thus, implementation for Case B is always better than implementation for Case A.

C. Case C—Preemptible Background Processes

The state transition diagram is shown in Fig. 4. Here the client is blocked for significantly shorter durations (more details in Section V). In this implementation, the foreground process can interrupt the background processes. The background processes then continue to do the actual updates in states $s_{1,2}, s_{2,2}, \dots, s_{1,2}, \dots$. The

I



Fig. 4. State transition for Case C.

general equations are

for $s_{i+1,1}$, $i \ge 0$: $\mu_i P_{i+1,1} = \lambda P_{i,2}$ for $s_{i,2}$, $i \ge 1$: $(\lambda + \mu'_i) P_{i,2} = p' \mu_{i-1} P_{i,1} + p \mu_i P_{i+1,1}$ $\sum_{i=0}^{\infty} P_{i,2} + \sum_{i=1}^{\infty} P_{i,1} = 1.$

And these can be solved to give

$$P_{0,2} = \frac{1}{\sum_{i=0}^{\infty} (1 + \lambda/\mu_i) \prod_{j=1}^{i} \left[p'\lambda/(p'\lambda + \mu'_j) \right]}$$

And we have

$$PW = \sum_{i=0}^{\infty} P_{i,2} = \left[\sum_{i=0}^{\infty} \prod_{j=1}^{i} \left[p'\lambda/(p'\lambda + \mu_j')\right]\right] P_{0,2}$$

Now

$$R = \frac{\sum_{i=0}^{\infty} (1/\mu_i) P_{i,2}}{\sum_{i=0}^{\infty} P_{i,2}} = \frac{\sum_{i=0}^{\infty} (1/\mu_i) \prod_{j=1}^{i} [p'\lambda/(p'\lambda + \mu'_j)]}{\sum_{i=0}^{\infty} \prod_{j=1}^{i} [p'\lambda/(p'\lambda + \mu'_j)]}$$

A closed form solution is useful to draw some conclusions regarding the performance of this implementation. This is possible under the following assumptions.

1) $\mu'_i = \mu = \text{constant}.$

2) $\mu_i = \mu_0/(i+1)$.

Assumption 1 is optimistic, although it is a reasonable approximation for certain update operations such as removing dead elements. Assumption 2 is pessimistic for a number of processor structures; a more reasonable assumption would be $\mu_i = [\mu_{01}/\log (i + 1)] + \mu_{02}$ where μ_{02} is the average rate for searching the sorted array while μ_{01} is the average rate for searching the binary tree. We

also assume that a steady state is reached so that probabilities for the operations ADD and DELETE are equal; otherwise, the search table will overflow or be empty in the steady state.

With these assumptions we obtain

$$R = \frac{1}{\mu_0} \left(1 + \frac{p'\lambda}{\mu} \right)$$
$$PW = \frac{1/\lambda}{1/\lambda + (1 + p'\lambda/\mu)/\mu_0}.$$

For small λ , these can be simplified to give

$$R \approx \frac{1}{\mu_0}$$
$$PW \approx \frac{1/\lambda}{1/\lambda + 1/\mu_0}$$

Since this is the best possible response time, this implementation is very efficient for reasonable rates of client requests. However, as λ increases, *R* tends to infinity and *PW* tends to 0. Thus, at some point the performance of Case C becomes worse than that of Case A. Notice that the three models presented here provide varying levels of concurrency by varying the degree of atomicity.

V. PROGRAM FOR CASE C

The implementation of Case A is straightforward and the implementation of Case B is a special case of that of Case C. So, we now present the details of the implementation for Case C.

A. Data Structures

The data structure is given below and consists of two binary search trees (with 11, t2 pointing to the roots, and for conciseness we will call them tree t1 and tree t2, respectively) and two arrays, a1 and a2. The reason for using two binary search trees is as follows. The foreground process adds an element to tree t1 when executing an ADD operation. The job of the background process is to merge the contents of a tree with the elements of an array. However, since the background process is preemptible we should allow for elements to be added while this merging is in progress. Hence, the usage of two binary search trees.

Of the two arrays, we use al for copying the elements from tree t2 in inorder traversal. This array is not really necessary but its usage simplifies the code. The array a2 contains a sorted list of elements. This array actually has two partitions, a21 and a22, so that the background process alternates between merging a21 and a1 into a22, and merging a22 and a1 into a21. In this process of merging, a21 and a22 grow in opposite directions. We assume that sufficient space is allocated for a2 so that the two partitions do not overlap.

It is possible to write an in-place merge of two array wherein the elements of two sorted arrays are merged without making use of a third array. Such a scheme would

T

avoid extra storage and is particularly efficient for this problem as one of the arrays contains a number of "dead" elements and hence the data movement can be made minimal. We wrote such an in-place merge, however, for simplicity we did not use it in our simulation studies and hence we do not present the code here.

Data Structure:

start1, end1: 0.. num1;

var a2: array [1..num2] of element; start21,end21,start22,end22: integer

The amount of additional storage required for the short term storage depends on the arrival rate of the various requests.

B. Program

For simpler and cleaner presentation we take some liberties with Pascal syntax as summarized below.

1) "{ " and "}" are used for "begin" and "end", respectively; these will also be used for enclosing assertions. The usage will be clear from the context.

2) Comments are preceded by "-". Explicit return statement is used.

t1:= nil;	-initialization
t2:= nil;	
start1:= 1; end1:= 0;	
start21:= 1; end21:= 0;	-start21 is always 1
start22: = $num2 + 1$; end22: = $num2$	-end22 is always num2

Foreground Process:

LOOKUP(k): found X p:= search_data_structures(k,false); if not found or else p.S = dead then return error message else return p.I;

ADD(k,i): found X p:= search_data_structures(k,true); if found and p.S = alive then return error message else {p.I:= i; p.S:= alive}

DELETE(k): found X p:= search_data_structures(k,false);
 if not found or else p.S = dead then return error message
 else p.S:= dead;

search_data_structures(k: key; add: boolean):
found X p:= if end21 < start21 or a2[end21].K < k</pre>

3) A number of values can be returned as a result of a procedure call, for example, use the notation ret_value1 X ret_value2 := procedure_name(..).

4) To maintain uniformity between arrays and trees we use "address of a[m]" to refer to the index "m" in the array a. For simplicity, if p is an address we use p.I etc.

-Two partitions: a21 = a2[start21..end21]- a22 = a2[start22..end22]

to refer to elements in trees as well as elements in the arrays.

5) "<" and ">" are used for enclosing atomic actions.

The main features of the program presented below are the following. A search for an element proceeds by searching a2, (a1 or t2), t1, in that order until the search is successful. Further, the array a2 is searched by searching only one of its two partitions. The background process ensures that new elements are added only to t1 by switching around the pointers to the roots of the two trees.

```
then binary_search(a2,start22,end22)
                     else binary search(a2,start21,end21);
     if found and p.S = alive then return true X p;
     found X p: = if end1 < start1 or a1[end1].K < k
                     then search tree(t2,false)
                     else binary search(a1,start1,end1);
     if found and p.S = alive then return true X p;
     return search tree(t1,add);
   binary_search(var a: array [ < > ] of element; 1,u: integer):
      while l < = u do
         {m := (1 + u) \text{ div } 2;}
        if a[m].K = k then return true X address of a[m]
        else if a[m].K < k then l := m + 1
        else u := m - 1;
      return false X nil;
   search_tree(var t: ptr; add: boolean);
     if t = nil then
        if not add then return false X nil
        else {new(t); t^{.}E.K := k; return false X address of t^{.}E}
      else if t^{\cdot}.E.K = k then return true X address of t^{\cdot}.E
      else if t<sup>.</sup>.E.K < k then search_tree(t<sup>.</sup>.right,add)
      else search_tree(t^.left,add);
Background Process:
   while true do
      \{ < t2 := t1; t1 := nil; >
      transfer t2 to a1;
      expand a22;
      < t2: = t1; t1: = nil; >
      transfer_t2_to_a1;
      expand a21;
      end1:= 0; start1:= 1\};
   transfer t2 to a1: inorder(t2)
   inorder(var t: ptr) : if t < > nil then
                            {inorder(t^.left):
                            \langle if t^{.}E.S. = alive then \{end1 := end1 + 1; a1[end1] := t^{.}E\};
                            \mathbf{x} := \mathbf{t};
                           t := t^{,right} >;
                            dispose(x);
                            inorder(t)
   expand a22: while end_{21} > = start_{21} and end_{12} > = start_{12} do
                    if a2[end21].S = dead then end21:= end21-1
                    else if a1[end1].S = dead then end1:= end1-1
                    else if a2[end21].K > a1[end1].K
                       then \{ < \text{start22} := \text{start22} - 1; a2[\text{start22}] := a2[\text{end21}]; \text{end21} := \text{end21} - 1 > \}
                    else
                               -a2[end21].K < a1[end1].K
                      \{ < start22 := start22 - 1; a2[start22] := a1[end1]; end1 := end1 - 1 > \};
                  while end_{21} > = start_{21} do
                    if a2[end21].S = alive
                    then \{ < \text{start22} := \text{start22} - 1; a2[\text{start22}] := a2[\text{end21}]; \text{end21} := \text{end21} - 1 > \}
                    else end21:= end21-1;
                  while end 1 > = start 1 do
                    if al[end1].S = alive
                    then \{ < \text{start22} := \text{start22} - 1; \text{ a2[start22]} := \text{a1[end1]}; \text{end1} := \text{end1} - 1 > \}
                    else end1: = end1 - 1;
```

expand a21: while end22 > = start22 and end1 > = start1 do if a2[start22].S = dead then start22:= start22 + 1else if a1[start1].S = dead then start1:= start1 + 1 else if $a_{start22}$.K < a_{start1} .K then { < end21: = end21 + 1; a2[end21]: = a2[start22]; start22: = start22 + 1 > } -a2[start22].K < a1[start1].Kelse $\{ < end_{21} := end_{21} + 1; a_{2}[end_{21}] := a_{1}[start_{1}]; start_{1} := start_{1} + 1 > \};$ while $end_{22} > = start_{22} do$ if a2[end22].S = alive then { < end21: = end21 + 1; a2[end21]: = a2[start22]; start22: = start22 + 1 > } else start22:= start22 + 1; while end 1 > = start 1 do if al[start1].S = alivethen $\{ < \text{end}_{21} := \text{end}_{21} + 1; a_2[\text{end}_{21}] := a_1[\text{start}_1]; \text{start}_1 := \text{start}_1 + 1 > \}$ else start1: = start1 + 1; VI. PROOF OF CORRECTNESS From INV it follows that to search for an element in

We discuss the correctness of the program presented in the previous section. To avoid excessive details we will only sketch the correctness proof leaving out some of the details. Let From INV it follows that to search for an element in a2[start21..end21] and a2[start22..end22], it is enough to search only one partition. Similarly, it is enough to search for an element in a1[start1..end1] and t2 by just searching one of them. Proving A1 is then straightforward. In prov-

D = The multiset of all the alive elements in a1[start1..end1], a2[start21..end21], a2[start22..end22], tree t1 and tree t2.
 INV = (No duplicates in D) ∧ (D is the set of all the elements in the system)

- ∧ (a1[start1..end1] ∘ inorder(t2), a2[start21..end21] ∘ a2[start22.end22], inorder(t1) are all sorted on key)
 - \wedge start21=1 \wedge end22=num2

where "o" denotes concatenation.

The invariant INV is actually the weak invariant that we referred to earlier. The strong invariant for this program, SINV, is given below ing A2 the only complication is that if a new location is allocated we require that it be in tree t1 and this is achieved because search of t2 is done always by search_tree(t2, false).

 $\begin{aligned} \text{SINV} &= (\text{No duplicates in D}) \land (\text{D is the set of all the elements in the system}) \\ \land (\text{there are no dead elements in the system}) \\ \land \text{t1} = \text{nil} \land \text{t2} = \text{nil} \land \text{end1} < \text{start1} \land [\text{end21} < \text{start21} \lor \text{end22} < \text{start22}] \\ \land (a2[\text{start21..end21}] \circ a2[\text{start22..end22}] \text{ is sorted on key}). \end{aligned}$

The background process in Case B establishes the strong invariant; however, in Case C this is established only if the background process is not preempted by the foreground process for a period of time. We will come back to this point later.

Consider the foreground process which is nonpreemptible. We first prove that INV is indeed an invariant. For this we first show that Now, LOOKUP and DELETE both maintain the invariant INV since search_data_structures(k,false) maintains it. For ADD, executing "p.I:= i; p.S:= alive" maintains the invariant INV because of A2.

We now consider the background process. Program annotation for the background process will be more complicated as this process can be preempted by the foreground process. We show the following annotation to be valid.

A1 : {INV} found X p := search_data_structures(k,false) {INV \land (found \Rightarrow p.K = k) \land (\neg found \Rightarrow (\forall j \in D : j.K \neq k))} A2 : {INV} found X p := search_data_structures(k,true) {INV \land (found \Rightarrow p.K = k) \land (\neg found \Rightarrow (\forall j \in D - {p}: j.K \neq k \land (p points to a new location added to t1)))}

T

11 : {INV \wedge t2 = nil \wedge start1=1 \wedge end1=0 \wedge start22=num2+1} 12 : while true do 13 • {INV \land t2 = nil \land start1 = 1 \land end1 = 0 \land start22 = num2 + 1} $\{ < t2 := t1; t1 := nil; >$ 14 • 15 {INV \land start1=1 \land end1=0} : 16 : transfer t2_to_a1; 17 ${INV \land t2 = nil \land start1 = 1}$: 18 expand_a22; ٠ 19 : {INV \wedge t2 = nil \wedge start1 = 1 \wedge end1 = 0 \wedge end21 = 0} $\{ < t2 := t1; t1 := nil; >$ 110:{INV \land start1=1 \land end1=0} 111 : 112 : transfer_t2_to_a1; 113 : ${INV \land t2 = nil \land start1 = 1}$ 114 : expand a21; 115 : {INV \wedge t2 = nil \wedge end1 < start1 \wedge start22=num2+ 1} 116: end1:=0; start1:=1{INV \land t2 = nil \land start1 = 1 \land end1 = 0 \land start22 = num2 + 1} 117 : };

We now proceed sequentially to show that the above annotation is valid. Assertion at 11 holds because of the initialization. We now show that

$$INV \wedge t2 = nil \wedge start1 = 1 \wedge end1 = 0 \wedge start22$$
$$= num2 + 1$$

is the loop invariant for the background process. We do that by assuming it holds at the start of the execution of this loop and show that it will then hold after the loop has been executed once more. The validity of the triples 13,14,15; 19,110,111, and 115,116,117 is straightforward. We therefore consider the assertions for transfer_t2_to_a1, expand_a22 and expand_a21.

The effect of the code for transfer_t2_to_a1 is that values from t2 are copied in inorder traversal into a1. Now since the precondition for executing this is that start1=1 and end1=0 it follows that $a1[start1..end1] \circ t2$ is sorted; everything in INV remains as before.

For the code expand_a22, the following loop invariants can be used for the three **while** loops respectively.

Consequently the background process maintains the invariant INV.

We now prove that the foreground and background processes correctly implement their respective tasks.

Foreground Process: Its "rely condition" is that the multiset D is not changed by the background process. The background process "guarantees" this since 1) whenever it executes "<t2:=t1; t1:=nil>," t2 is nil, and 2) all data structure modifications required in order to transfer items with flag=alive from one data structure to another (e.g., in procedures transfer_t2_to_a1, expand_a22, expand_a11) are atomic and correct.

Background Process: Its "rely condition" is that the only modifications by the foreground process to t2, a1, a21, a22 is to change the flag of an item from alive to dead. The foreground process "guarantees" this since the only modifications to t2, a1, a21, a22 is to change the flag of an item from alive to dead in operation DELETE(k). All the other modifications by the foreground process oc-

- I1: (a2[start21..end21] ∘ a2[start22..end22], a1[start1..end1] ∘ a2[start22..end22] is sorted on key) ∧ t2=nil
- I2: I1 \land ((start1=1 \land end1=0) \lor end21=0) I3: I1 \land end21=0

and consequently the following assertion will hold at the end of expand a22

INV \wedge t2=nil \wedge start1=1 \wedge end1=0 \wedge end21=0

For the code expand_a21, the following loop invariants can be used for the three **while** loops, respectively.

cur in t1, and these are of no consequence to the background process.

The task of the background process is to achieve the strong invariant, SINV, if it is not interrupted by the foreground process. Assume that it is not interrupted between 13 and 19. Then, 1) t1 is nil since the foreground process

I4: (a2[start21..end21] \circ a2[start22..end22], a2[start21..end21] \circ a1[start1..end1] is sorted on key) \wedge t2=nil

I5: I4 \land (end1 < start1 \lor start22=num2+ 1) I6: I4 \land start22=num2+ 1

and consequently the following assertion will hold at the end of expand a21

 $INV \wedge t2 = nil \wedge end1 < start1 \wedge start22 = num2 + 1$

has not added any items to t1 (by assumption), and 2) there are no items with flag=dead since a) all such items which existed at 13 have been removed by the background



process by the time execution reaches 19, and b) the foreground process has not changed the flag of any item (by assumption). Hence, at 19 we can assert:

INV \wedge t2=nil \wedge start1=1 \wedge end1=0 \wedge end21=0 \wedge t1=nil \wedge (there are no items with flag=dead in the system)

This implies SINV. A similar argument can be given for the phase starting at 19 and terminating at 117.

VII. PERFORMANCE STUDIES

Fig. 5 shows the results of an experiment conducted to compare the performance of cases A and C using VAX[®]/VMS interrupt and timing system services. The details of the experiment are discussed in [12]. The parameters of the experiment are as follows.

1) λ , the rate at which the client issues a request.

2) p, the probability that the request is for LOOKUP.

Also, for each request there is a parameter which specifies the probability that the request is valid. For the data shown here, these probabilities are all 0.9.

From Fig. 5, we observe that

1) For small to medium λ , Case C is much better than Case A.

2) Case C is better than Case A as the LOOKUP probability p increases.

VIII. SUMMARY

In this paper we have discussed one method for developing high performance implementation for abstract data types. This method relies on the usage of multilevel data structures and maintenance processes and achieves high performance by scheduling the computations required to maintain the data structure during the idle periods. This approach is particularly suitable for distributed system as the data reorganization can be done while response is

*VAX is a registered trademark of Digital Equipment Corporation.

being sent to the client. We have presented various models for the maintenance processes, the different models provide varying amounts of concurrency. Both the derived performance and the experimental results favor the usage of maintenance processes with finer grained concurrency over coarse grained concurrency for a wide range of parameters.

Possible research directions in this area include developing performance models for multilevel data structures when the foreground process is also preemptible and when there are multiple clients concurrently accessing the data type.

REFERENCES

- F. B. Bastani, S. S. Iyengar, and I. L. Yen, "Concurrent maintenance of multilevel data structures," Dep. Comput. Sci., Univ. Houston-University Park, Tech. Rep. UH-CS-85-3, Jan. 1985.
- [2] F. B. Bastani, I. L. Yen, A. Moitra, and S. S. Iyengar, "Impact of parallel processing on software quality," in *Proc. First Int. Conf. Supercomputing Systems*, Dec. 1985, pp. 369-376.
- [3] E. W. Dijkstra, L. Lamport, A. J. Martin, C. S. Sholten, and E. F. M. Steffens, "On-the-fly garbage collection: An exercise in cooperation," Commun. ACM, vol. 21, no. 11, pp. 966-975, Nov. 1978.
- [4] C. S. Ellis, "Distributed data structures: A case study," in Proc. 5th Int. Conf. Distributed Processing Systems, May 1985, pp. 201-209.
- [5] T. Hickey and J. Cohen, "Performance analysis of on-the-fly garbage collection," *Commun. ACM*, vol. 27, no. 11, pp. 1143–1154, Nov. 1984.
- [6] C. B. Jones, "Tentative steps towards a development method for interfering programs," ACM Trans. Program. Lang. and Syst., vol. 5, no. 4, pp. 596-619, Oct. 1983.
- [7] L. Kleinrock, Queueing Systems-Vol. 1: Theory. New York: Wiley, 1975.
- [8] B. W. Lampson, "Hints for computer system design," IEEE Software, vol. 1, no. 1, pp. 11-28, Jan. 1984.
- [9] U. Manber, "Concurrent maintenance of binary search trees," *IEEE Trans. Software Eng.*, vol. SE-10, no. 6, pp. 777-784, Nov. 1984.
- [10] U. Manber and R. E. Ladner, "Concurrency control in a dynamic search structure," ACM Trans. Database Syst., vol. 9, no. 3, pp. 439-455, Sept. 1984.
- [11] S. Owicki and D. Gries, "Verifying properties of parallel programs: An axiomatic approach," *Commun. ACM*, vol. 19, no. 5, pp. 279-285, May 1976.
- [12] I. L. Yen, "The role of parallel processing in application programs," M.S. thesis, Dep. Comput. Sci., Univ. Houston-University Park, May 1985.

I



Abha Moitra received the M.Sc. degree in physics from the Birla Institute of Technology and Science in 1977, and the Ph.D. degree in computer science from the University of Bombay, Bombay, India, in 1981.

She is currently an Assistant Professor of Computer Science at Cornell University, Ithaca, NY. She has authored several papers in distributed computing, parallel algorithms, and other related areas of computer science.



S. Sitharama Iyengar received the Ph.D. degree in engineering in 1974.

He is currently a Professor of Computer Science and Supervisor of robotic research and parallel algorithms at Louisiana State University. He has authored (or coauthored) more than 70 research papers in parallel algorithms, data structures, and navigation of intelligent mobile robots. He is currently studying the application of neural network techniques for path planning and learning in mobile robots. His papers have appeared in the

following journals: IEEE-TSE, IEEE-PAMI, IEEE-SMC, IEEE JOURNAL OF ROBOTICS AND AUTOMATION. IEEE TRANSACTIONS ON COMPUTERS, CACM, JCIS, Computer Networks, Journal of Robotic Systems, BIT, Theoretical Computer Science, and several other international journals and IEEE proceedings. His research has been funded by NASA, DOE, the Navy, etc.

Dr. Iyengar is an ACM National Lecturer for 1986-1988. He has been on the program committee for several major conferences in the U.S. and in Europe.

Farokh B. Bastani (M'82) received the B.Tech. degree in electrical engineering from the Indian Institute of Technology, Bombay, India, in 1977, and the M.S. and Ph.D. degrees in electrical engineering and computer science from the University of California, Berkeley, in 1978 and 1980, respectively.

He joined the University of Houston, Houston, TX, in 1980, where he is currently Assistant Professor of Computer Science. His research interests include developing a design methodology and quality assessment techniques for large scale computer systems.

I. Ling Yen received the B.S. degree in physics from National Tsing-Hua University, Republic of China, and the M.S. degree in computer science from the University of Houston, Houston, TX.

She is currently a Systems Programmer at the University of Massachusetts at Amherst. Her current interests are in distributed computing and artificial intelligence.

Ms. Yen is a member of the Association for Computing Machinery.