

1 Robust BST

The original BST algorithm does not support the crash-recovery model. It is clear from the code a process does not persist the operation's response in the non-volatile memory, and thus, once a process crash the response is lost. For example, assume a process q apply $\text{INSERT}(k)$, performs a successful CAS in line 96 and fails after completing the HELPINSERT routine. In this case, the INSERT operation took effect, that is, the new key appears as a leaf in the tree, and any $\text{FIND}(k)$ operation will return it. However, even though the operation must be linearized before the crash, upon recovery process q is unaware of it. Moreover, looking for the new leaf in the tree may be futile, as it might be k has been removed from the tree after the crash.

Furthermore, if no recover routine is supplied, it may result an execution which is not well-formed. Consider for example the following scenario. A process q invoke an $Op_1 = \text{INSERT}(k_1)$ operation. q performs a successful CAS in line 96 followed by a crash. After recovering, q invoke an $Op_2 = \text{INSERT}(k_2)$ operation. Assume k_1 and k_2 belongs to a different parts of the tree (do not share parent or grandparent). Then, q can complete the insertion of k_2 without having any affect on k_1 . Now, a process q' performs $\text{FIND}(k_1)$ which returns \perp , as the insertion of k_1 is not completed, followed by $\text{FIND}(k_2)$, which returns the leaf of k_2 . The $\text{INSERT}(k_1)$ operation will be completed later by any INSERT or DELETE operation which needs to make changes to the flagged node. We get that Op_2 must be before Op_1 in the linearization, although Op_1 invoked first.

The kind of anomaly described above can be addressed by having the first CAS of a successful attempt for INSERT or DELETE as the linearization point, as in the Linked-List. For that, the FIND routine should take into consideration future unavoidable changes, for example, a node flagged with IFlag ensures an insertion of some key. A simple solution is to change the FIND routine such that it also helps other operations, as described in figure 1. The FIND routine will search for key k in the tree. If the SEARCH routine returns a grandparent or a parent that is flagged, then it might be that an insert or delete of k is currently in progress, thus we first help the operation to complete, and then search for k again. Otherwise, if $gpupdate$ or $pupdate$ has been changed since the last read, it means some change already took affect, and there is a need to search for k again. If none of the above holds, there is a point in time where gp points to p which points to l , and there is no attempt to change this part of the tree. As a result, if k is in the tree at this point, it must be in l , and the find can return safely.

The approach described above is not efficient in terms of time. We would like a solution which maintain the desirable behaviour of the original FIND routine, where a single SEARCH is needed. A more refined solution is given in figure 2. The intuition for it is drawn from the Linked-List algorithm. In the Linked-List algorithm it was enough to consider a marked node as if it has been deleted, without the need to complete the deletion. Nonetheless, the complex BST implementation is more challenging, as the DELETE routine needs to successfully capture two nodes using CAS in order to complete the deletion. Therefore, if a process p executes $\text{FIND}(k)$ procedure, and observes a node flagged with DFlag attempting to delete the key k , it can not know whether in the future this delete attempt will succeed or fail, and thus does not know whether to consider the key k as part of the tree or not. To overcome this problem, in such case the process will first try and validate the delete operation by marking the relevant node. According to whether the marking attempt was successful, the process can conclude if the delete operation is successful or not. In order to easily implement the modified FIND routine there is a need to conclude from IInfo what is the new leaf (leaf *new* in the INSERT routine). For simplicity of presentation, we do not add this field, and abstractly refer to it in the code.

The correctness of the two suggested solutions relies on the following argument. Once a process flags a node during operation Op with input key k (either INSERT or DELETE), then if this attempt to complete the operation eventually succeed (i.e., the marking is also successful in the case of DELETE), then any FIND(k) operation invoked from this point consider Op as if it is completed.

The suggested modification, although being simple and local, only guarantee the implementation satisfy R-linearizability. However, the problem of response being lost in case of a crash is not addressed. Roughly speaking, the critical points in the code for recovery are the CAS primitives, as a crash right after applying CAS operation results the lost of the response, and in order to complete the operation the process needs to know the result of the CAS. In addition, because of the helping mechanism, a suspended DELETE operation which flagged a node and yet to mark one, may be completed by other process in the future, and may not. Upon recovery, the process needs to distinguish between the two cases, in order to obtain the right response.

To address this issue, we expend the helping mechanism, so that the helping process needs also to update the info structure in case of a success. This is done by adding a boolean field to the Flag structure. This way, if a process crash along an operation Op , upon recovery it can check whether the operation was completed by some different process.

Before a process attempt to perform an operation, as it creates the Flag structure op describing the operation and its affect on the data structure, the process stores op in a designated location (for simplicity, we use an array). Upon recovery, the process reads this location, and if the operation is not completed yet, it retries to perform it, starting from the point of the first flagging (the first CAS). Otherwise, the operation was completed, and the response value is already known. Notice that there is a scenario in which process q recovers and observes an operation Op as not being completed, but just before it retries it, some other process complete the operation. We need to prove that even in such case, the operation will affect the data structure exactly once, and the right response is returned.

Notice that the given implementation does not recover the FIND routine, since this routine does not make any changes to the BST, hence it is always safe to consider it as having no linearization point and reissue it. Also, for ease of presentation, we only write to Announce[id] once we are about to capture a node using a CAS. However, writing to Announce[id] at the beginning of the routine may be helpful in case of a crash early in the routine, so that the process will be able to use the data stored in Announce[id] in such case also. The same is true with response value - Announce[id].done is updated only if the routine made changes to the BST.

1.0.1 Correctness

In this subsection we give a proof sketch for the algorithm correctness. The proof relies on the correctness of the original algorithm, which can be found on [...]. Moreover, the original algorithm is anonymous and uniform, i.e., any number of processes can use the BST, and there is no need for a process to know the number of processes in the system in order to use the BST. Thus, if we consider an execution and prove it is indistinguishable to a process p from some execution of the original algorithm in which more processes can participate, then it has to return the same response on both.

It is easy to verify the post-conditions of the SEARCH routine still holds, as they follow directly from the code. Also, since no changes are made to the SEARCH routine, it does not make changes to the BST, but simply traverse it. Therefore FIND routine, which only uses SEARCH, does not affect any process, and in case of a failure along FIND execution, reissuing it satisfies NRL. A key

argument in the proof is that a node in the tree never point twice to the same node. This attributes to the fact an INSERT presents a node with two fresh children. Therefore, if we assume node nd_1 points to nd_2 , in order to point to a different node we must either delete nd_2 , or insert a new node between the two. In the last case, a copy of nd_2 is created, and thus nd_2 will no longer be used.

Consider a process q performs an INSERT(k) operation. As argued before, a crash before writing to $\text{Announce}[q]$ implies no changes has been made to the BST. Assume thus q executes line 95, i.e., q stored in $\text{Announce}[q]$ a pointer to an IInfo record containing all the data needed for the current attempt to complete the INSERT routine. Since q is in the middle of a while loop, it is enough to prove that if q crash before the next time it writes to $\text{Announce}[q]$, if there is such write, upon recovery it will either complete its operation with the right response, or will continue to the next write to $\text{Announce}[q]$ without having any affect on the BST. Hence, the same argument can be applied once q writes to $\text{Announce}[q]$ again.

Assume q performs a successful CAS in line 96. Then, the IInfo op is stored in $p \leftarrow \text{update}$, and it is also iflagged. Notice that p was not changed since the SEARCH routine read it. This follows from the fact that any change to a node must store a pointer to a new created record, different then any other record. Thus, a successful CAS implies no changes have been applied. Starting from this point, no changes can be made to p , except for the change point to by op, as the node is flagged. Relying on the correctness of the original algorithm, no matter how many times $\text{HELPINSERT}(\text{op})$ will be executed, the change will occur only once. This follows from the fact that many process can observe op, and will try to complete it in the future. The core for this argument is that a node never point twice to the same node.

```

FIND(Key  $k$ ) : Leaf* {
1      Internal * $gp$ , * $p$ 
2      Leaf * $l$ 
3      Update  $pupdate, gpupdate$ 
4      while (TRUE) {
5           $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
6          if  $gpupdate.state \neq \text{CLEAN}$  then  $\text{HELP}(gpupdate)$ 
7          else if  $pupdate.state \neq \text{CLEAN}$  then  $\text{HELP}(pupdate)$ 
8          else if  $gp \leftarrow update = gpupdate$  and  $p \leftarrow update = pupdate$  then {
9              if  $l \rightarrow key = k$  then return  $l$ 
10             else return  $\perp$ 
11         }
12     }
13 }
```

Figure 1: Solution 1: R-linearizable FIND routine

```

FIND(Key k) : Leaf* {
14   Internal *gp, *p
15   Leaf *l
16   Update pupdate, gpupdate

17    $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
18   if  $l \rightarrow key \neq k$  then {
19       if (pupdate.state = IFlag and pupdate.info attempt to add key k) then
20           return leaf with key k from pupdate.info
21       else return  $\perp$ 
22   }
23   if (pupdate.state = MARK and pupdate.info  $\leftarrow l \leftarrow key = k$ ) then return  $\perp$ 
24   if (gpupdate.state = DFlag and gpupdate.info  $\leftarrow l \leftarrow key = k$ ) then {
25       op := gpupdate.info
26       result := CAS(op  $\rightarrow p \rightarrow update$ , op  $\rightarrow pupdate$ ,  $\langle \text{MARK}, op \rangle$ )    ▷ mark CAS
27       if (result = op  $\rightarrow pupdate$  or result =  $\langle \text{MARK}, op \rangle$ ) then return  $\perp$     ▷ op  $\rightarrow p$  is successfully marked
28   }
29   return l
30 }

```

Figure 2: Solution 2: R-linearizable FIND routine

```

31 type Update {                ▷ stored in one CAS word
32     {CLEAN, DFlag, IFlag, MARK} state
33     Flag *info
34 }
35 type Internal {              ▷ subtype of Node
36      $Key \cup \{\infty_1, \infty_2\}$  key
37     Update update
38     Node *left, *right
39 }
40 type Leaf {                  ▷ subtype of Node
41      $Key \cup \{\infty_1, \infty_2\}$  key
42 }
43 type IInfo {                 ▷ subtype of Flag
44     Internal *p, *newInternal
45     Leaf *l
46     boolean done
47 }
48 type DInfo {                 ▷ subtype of Flag
49     Internal *gp, *p
50     Leaf *l
51     Update pupdate
52     boolean done
53 }
▷ Initialization:
54 shared Internal *Root := pointer to new Internal node
    with key field  $\infty_2$ , update field  $\langle \text{CLEAN}, \perp \rangle$ , and
    pointers to new Leaf nodes with keys  $\infty_1$  and
     $\infty_2$ , respectively, as left and right fields.

```

Figure 3: Type definitions and initialization.

```

RECOVER() {
55     Flag *op = Announce[id]

56     if op ← done = TRUE then return TRUE
57     if op of type IInfo then
58         go to line 96
59     if op of type DInfo then
60         go to line 125
61 }

```

Figure 4: RECOVER routine

```

62 SEARCH(Key k) : (Internal*, Internal*, Leaf*, Update, Update) {
    ▷ Used by INSERT, DELETE and FIND to traverse a branch of the BST; satisfies following postconditions:
    ▷ (1)  $l$  points to a Leaf node and  $p$  points to an Internal node
    ▷ (2) Either  $p \rightarrow \text{left}$  has contained  $l$  (if  $k < p \rightarrow \text{key}$ ) or  $p \rightarrow \text{right}$  has contained  $l$  (if  $k \geq p \rightarrow \text{key}$ )
    ▷ (3)  $p \rightarrow \text{update}$  has contained  $pupdate$ 
    ▷ (4) if  $l \rightarrow \text{key} \neq \infty_1$ , then the following three statements hold:
    ▷ (4a)  $gp$  points to an Internal node
    ▷ (4b) either  $gp \rightarrow \text{left}$  has contained  $p$  (if  $k < gp \rightarrow \text{key}$ ) or  $gp \rightarrow \text{right}$  has contained  $p$  (if  $k \geq gp \rightarrow \text{key}$ )
    ▷ (4c)  $gp \rightarrow \text{update}$  has contained  $gpupdate$ 
63   Internal *gp, *p
64   Node *l := Root
65   Update gpupdate, pupdate
66   while  $l$  points to an internal node {
67       gp := p
68       p := l
69       gpupdate := pupdate
70       pupdate := p → update
71       if  $k < l \rightarrow \text{key}$  then  $l := p \rightarrow \text{left}$  else  $l := p \rightarrow \text{right}$ 
72   }
73   return (gp, p, l, pupdate, gpupdate)
74 }
75 FIND(Key k) : Leaf* {
76   Leaf *l
77   (−, −, l, −, −) := SEARCH(k)
78   if  $l \rightarrow \text{key} = k$  then return  $l$ 
79   else return  $\perp$ 
80 }
81 INSERT(Key k) : boolean {
82   Internal *p, *newInternal
83   Leaf *l, *newSibling
84   Leaf *new := pointer to a new Leaf node whose key field is  $k$ 
85   Update pupdate, result
86   IInfo *op
87   while TRUE {
88       (−, p, l, pupdate, −) := SEARCH(k)
89       if  $l \rightarrow \text{key} = k$  then return FALSE
90       if pupdate.state ≠ CLEAN then HELP(pupdate)
91       else {
92           newSibling := pointer to a new Leaf whose key is  $l \rightarrow \text{key}$ 
93           newInternal := pointer to a new Internal node with key field  $\max(k, l \rightarrow \text{key})$ ,
94                       update field (CLEAN,  $\perp$ ), and with two child fields equal to new and newSibling
95                       (the one with the smaller key is the left child)
96           op := pointer to a new IInfo record containing (p, l, newInternal, FALSE)
97           Announce[id] := op
98           result := CAS(p → update, pupdate, (IFlag, op))
99           if result = pupdate or result = (IFlag, op) then {
100               HELPINSERT(op)
101               return TRUE
102           }
103           else HELP(result)
104       }
105   }
106 }
107 HELPINSERT(IInfo *op) {
108   ▷ Precondition:  $op$  points to an IInfo record (i.e., it is not  $\perp$ )
109   CAS-CHILD(op → p, op → l, op → newInternal)
110   CAS(op → p → update, (IFlag, op), (CLEAN, op))
111   op → done := TRUE
112 }

```

▷ Each stores a copy of an *update* field
 ▷ Remember parent of p
 ▷ Remember parent of l
 ▷ Remember *update* field of gp
 ▷ Remember *update* field of p
 ▷ Move down to appropriate child
 ▷ Cannot insert duplicate key
 ▷ Help the other operation
 ▷ iflag CAS
 ▷ The iflag CAS was successful
 ▷ Finish the insertion
 ▷ The iflag CAS failed; help the operation that caused failure
 ▷ ichild CAS
 ▷ iunflag CAS
 ▷ mark the operation as completed

Figure 5: Pseudocode for SEARCH, FIND and INSERT.

```

112 DELETE(Key  $k$ ) : boolean {
113   Internal * $gp$ , * $p$ 
114   Leaf * $l$ 
115   Update  $pupdate$ ,  $gpupdate$ ,  $result$ 
116   DInfo * $op$ 

117   while TRUE {
118      $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
119     if  $l \rightarrow key \neq k$  then return FALSE ▷ Key  $k$  is not in the tree
120     if  $gpupdate.state \neq \text{CLEAN}$  then HELP( $gpupdate$ )
121     else if  $pupdate.state \neq \text{CLEAN}$  then HELP( $pupdate$ )
122     else { ▷ Try to flag  $gp$ 
123        $op :=$  pointer to a new DInfo record containing  $\langle gp, p, l, pupdate, \text{FALSE} \rangle$ 
124        $\text{Announce}[id] := op$ 
125        $result := \text{CAS}(gp \rightarrow update, gpupdate, \langle \text{DFlag}, op \rangle)$  ▷ dflag CAS
126       if  $result = gpupdate$  or  $result = \langle \text{DFlag}, op \rangle$  then { ▷ CAS successful
127         if HELPDELETE( $op$ ) then return TRUE ▷ Either finish deletion or unflag
128       }
129       else HELP( $result$ ) ▷ The dflag CAS failed; help the operation that caused the failure
130     }
131     if  $op \rightarrow done = \text{TRUE}$  then
132       return TRUE
133   }
134 }

135 HELPDELETE(DInfo * $op$ ) : boolean {
136   ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
137   Update  $result$  ▷ Stores result of mark CAS
138    $result := \text{CAS}(op \rightarrow p \rightarrow update, op \rightarrow pupdate, \langle \text{MARK}, op \rangle)$  ▷ mark CAS
139   if  $result = op \rightarrow pupdate$  or  $result = \langle \text{MARK}, op \rangle$  then { ▷  $op \rightarrow p$  is successfully marked
140     HELPMARKED( $op$ ) ▷ Complete the deletion
141     return TRUE ▷ Tell DELETE routine it is done
142   }
143   else { ▷ The mark CAS failed
144     HELP( $result$ ) ▷ Help operation that caused failure
145      $\text{CAS}(op \rightarrow gp \rightarrow update, \langle \text{DFlag}, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ backtrack CAS
146     return FALSE ▷ Tell DELETE routine to try again
147   }

148 HELPMARKED(DInfo * $op$ ) {
149   ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
150   Node * $other$ 
151   ▷ Set  $other$  to point to the sibling of the node to which  $op \rightarrow l$  points
152   if  $op \rightarrow p \rightarrow right = op \rightarrow l$  then  $other := op \rightarrow p \rightarrow left$  else  $other := op \rightarrow p \rightarrow right$ 
153   ▷ Splice the node to which  $op \rightarrow p$  points out of the tree, replacing it by  $other$ 
154    $\text{CAS-CHILD}(op \rightarrow gp, op \rightarrow p, other)$  ▷ dchild CAS
155    $\text{CAS}(op \rightarrow gp \rightarrow update, \langle \text{DFlag}, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ dunflag CAS
156    $op \rightarrow done := \text{TRUE}$  ▷ mark the operation as completed
157 }

158 HELP(Update  $u$ ) { ▷ General-purpose helping routine
159   ▷ Precondition:  $u$  has been stored in the  $update$  field of some internal node
160   if  $u.state = \text{IFlag}$  then HELPINSERT( $u.info$ )
161   else if  $u.state = \text{MARK}$  then HELPMARKED( $u.info$ )
162   else if  $u.state = \text{DFlag}$  then HELPDELETE( $u.info$ )
163 }

164 CAS-CHILD(Internal * $parent$ , Node * $old$ , Node * $new$ ) {
165   ▷ Precondition:  $parent$  points to an Internal node and  $new$  points to a Node (i.e., neither is  $\perp$ )
166   ▷ This routine tries to change one of the child fields of the node that  $parent$  points to from  $old$  to  $new$ .
167   if  $new \rightarrow key < parent \rightarrow key$  then
168      $\text{CAS}(parent \rightarrow left, old, new)$ 
169   else
170      $\text{CAS}(parent \rightarrow right, old, new)$ 
171 }

```

Figure 6: Pseudocode for DELETE and some auxiliary routines.