

1 Linked-List

The original algorithm by Harris is presented in Figure 1. Harris approach uses an Atomic-Markable-Reference object, in which the next field of a Node, in addition to a reference to the next node in the list, is also marked or unmarked. The two fields can be update atomically, either together or individually. This can be done by using the most-significant-bit of next for the marking. For simplicity, we assume `node.next` returns the reference, while a query can be used to identify if it is marked. Therefore, whenever writing to `node.next`, or performing CAS, both the reference and marking state should be mention. For ease of presentation, we assume a List is initialised with head and tail, containing keys $-\infty, \infty$ respectively. We allow no insert or delete of these keys.

The Lookup procedure is used by Insert and Delete in order to find the node with the lowest key greater or equal to the input key, and its predecessor on the list, while physically removing any marked node on its way. To insert a key α , a process first finds the right location for α using the Lookup procedure, and then tries to set `pred.next` to point to a new node containing α by performing CAS. To delete a key α , a process looks for it using the Lookup procedure, and then tries to logically remove it by marking `curr.next` using CAS. In case the marking was successful, the process also tries to physically remove the node. To find a key α , a process simply looks for a node in the list with key α which is unmarked.

Procedure Lookup(int key)

Data: Node* pred, curr, succ

```
1 retry: while true do
2   | pred = head
3   | curr = head.next
4   | while true do
5   |   | succ = curr.next
6   |   | if curr.next is marked then
7   |   |   | if pred.next.CAS (unmarked curr, unmarked succ) == false then
8   |   |   |   | go to retry
9   |   |   | end
10  |   |   | curr = succ
11  |   | else
12  |   |   | if curr.key ≥ key then
13  |   |   |   | return ⟨pred, curr⟩
14  |   |   | end
15  |   |   | pred = curr
16  |   |   | curr = succ
17  |   | end
18  | end
19 end
```

Shared variables: Node* head

Procedure Insert(int key)	
<hr/>	
Data: Node* pred, curr Node node = new Node (key)	
20	while true do
21	⟨pred, curr⟩ = lookup(key)
22	if curr.key == key then
23	return false
24	else
25	node.next = unmarked curr
26	if pred.next.CAS (unmarked curr, unmarked node) then
27	return true
28	end
29	end
30	end
<hr/>	
Procedure Delete(int key)	
<hr/>	
Data: Node* pred, curr, succ	
31	while true do
32	⟨pred, curr⟩ = lookup(key)
33	if curr.key != key then
34	return false
35	else
36	succ = curr.next
37	if curr.next.CAS (unmarked succ, marked succ) then
38	pred.next.CAS (unmarked curr, unmarked succ)
39	return true
40	end
41	end
42	end
<hr/>	
Procedure Find(int key)	
<hr/>	
Data: Node* curr = head	
43	while curr.key < key do
44	curr = curr.next
45	end
46	return (curr.key == key && curr.next is unmarked)

Figure 1: Harris Non-Blocking Algorithm

1.0.1 Crash-Recovery

The linearization point are as follows:

Insert: At the point of a successful CAS

Delete: At the point of a successful CAS for marking the node (logical delete)

Find: At the point of the procedure return, that is, either when $\text{curr.key} \neq \text{key}$, or at the second condition test.

Following these linearization points (committing proof...), insert and delete operation are linearized at the point where they affect the system. That is, if there is a linearization point for insert operation, then all process will see the new node starting from this point, and if a node was logically removed, then all processes treat it as a removed node. Therefore, once a process p recovers following a crash, the List data structure is consistent - if p has a pending operation, either the operation already had a linearization point and affect all other processes, or it did not affect the data structure at all.

However, even though the List data structure is consistent, the response of the pending operation is lost. Consider for example a scenario in which process p performs $Delete(\alpha)$ and crash at line 37 after performing a successful CAS. Upon recovery, p may be able to decide α was removed, as the node is marked. Nevertheless, even if no other process takes steps, p is not able to determine whether it is the process to successfully delete α , or that it was done by some other process, and therefore it does not able to determine the right response. Moreover, in case the node was physically removed, p is not able to determine whether α has been deleted, as it is no longer part of the list.

1.0.2 Linked-List Recoverable Version

To solve the problems mention above, we present a modification for the algorithm such that in case a process fails, upon recovery it is able to complete its last pending operation and also return the response value.

Each node is equipped with a new field named deleter. This field is used to determine which process is the one to delete the node. After a process p successfully mark a node (logical delete), it tries to write its id to deleter using CAS. This way, if a process fails during a delete, it can use deleter in order to determine the response value. We assume deleter is initialized to null when creating a new node.

Each process p has a designated location in the memory, Backup[p]. Before trying to apply an operation, p writes to Backup[p] the entire data needed to complete the operation. Upon recovery, p can read Backup[p], and based on it to complete its pending operation, in case there is such. Formally, Backup[p] contains a pointer to a structure containing all the relevant data.

We present the modified algorithm. Only the procedures which require changes are presented. For simplicity, a process creates a new operation structure each time it writes to Backup, although a process can use two such structures alternately.

Correctness Argument

In the following, we give an intuition for the correctness of the algorithm.

For the Insert operation, p tries to add the new node by performing a CAS. If it succeeds it will return true if it suffers no failure. In case of a failure after writing to Backup[p], upon recovery p tries to complete its operation. If it already performed a successful CAS, that is, the node was added to the list, then either it is still in the list or that it was deleted. Therefore, if p can find the new node in the list (using a procedure similar to find), or that it is marked, it must be that the node was added, and p can return true. Otherwise, p either crashed before performing the CAS,

or that the CAS was unsuccessful. In both cases, the new node was not added, and p can restart the Insert procedure.

For the delete operation, once p logically delete a node v , it also tries to announce itself as the "removal" of the node by writing its id to deleter using CAS. Assume a process p crash while trying to delete the node. Upon recovery, if p sees the node is not marked, then obviously its deletion did not took affect, and it can restart the delete operation. However, if the node is marked, it might be that p marked it before the crash, or it might be some other process trying to delete the same node did so. As p can not distinguish between the two, and since we desire for a lock-free implementation, we let p to try and complete the deletion, even if it is not to process to logically delete the node. To avoid a scenario in which more then a single process "delete" the same node, they all compete for deleter using CAS. The first one to perform it will win, and it is the only process to return true. It is easy to verify once a process writes to deleter, then eventually, if given enough time with no crash, it returns true, while any other process trying to delete the same node will have to retry the delete operation.

```

Procedure Recover
73 if Backup[p].type == Insert then
74   if Backup[p].new is in the list || Backup[p].curr.next is marked then
75     return true
76   else
77     go to 47                                // restart Insert
78   end
79 end
80 end
81 if Backup[p].type == Delete then
82   if Backup[p].curr.next is marked then
83     go to 67                                // try to complete the deletation
84   else
85     go to 59                                // restart Delete
86   end
87 end
88 end
89 end

```

2 Elimination Stack

Procedure Exchange($\langle T, \text{flag} \rangle$ slot, T myitem, long timeout)

Data: long timeBound = getNanos() + timeout

```
90 while true do
91   if getNanos() > timeBound then
92     return TIMEOUT
93   end
94    $\langle \text{youritem}, \text{state} \rangle = \text{slot}$ 
95   switch state do
96     case EMPTY do
97       if slot.CAS ( $\langle \text{youritem}, \text{EMPTY} \rangle$ ,  $\langle \text{myitem}, \text{WAITING} \rangle$ ) then
98         while getNanos() < timeBound do
99            $\langle \text{youritem}, \text{state} \rangle = \text{slot}$ 
100          if state == BUSY then
101            slot =  $\langle \text{null}, \text{EMPTY} \rangle$ 
102            return youritem
103          end
104        end
105        if slot.CAS ( $\langle \text{myitem}, \text{WAITING} \rangle$ ,  $\langle \text{null}, \text{EMPTY} \rangle$ ) then
106          return TIMEOUT
107        else
108           $\langle \text{youritem}, \text{state} \rangle = \text{slot}$ 
109          slot =  $\langle \text{null}, \text{EMPTY} \rangle$ 
110          return youritem
111        end
112      end
113      break
114    end
115    case WAITING do
116      if slot.CAS ( $\langle \text{youritem}, \text{state} \rangle$ ,  $\langle \text{myitem}, \text{BUSY} \rangle$ ) then
117        return youritem
118      end
119      break
120    end
121    case BUSY do
122      break
123    end
124  end
125 end
```

Procedure TryPush(Node newNode)

```
126 Node* oldTop = Top
127 newNode.next = oldTop
128 if Top.CAS (oldTop, newNode) then
129   | return true
130 else
131   | return false
132 end
```

Procedure TryPop(void)

```
133 Node *oldTop = Top
134 Node *newTop
135 if oldTop == NULL then
136   | return EMPTY
137 end
138 newTop = oldTop.next
139 if Top.CAS (oldTop, newTop) then
140   | return oldTop
141 end
142 return NULL
```

3 BST

The original BST algorithm does not support the crash-recovery model. It is clear from the code a process does not persist an operation response in the non-volatile memory, and thus, once a process crash the response is lost. For example, assume a process p apply $\text{INSERT}(k)$ and assume p performs

Procedure push(T myitem)

```
143 Node *nd = new Node (myitem)
144 while true do
145   | if TryPush(nd) then
146     | return true
147   end
148   othervalue = visit(nd)
149   if othervalue == NULL then
150     | Record Success ()
151     | return true
152   else
153     | Record Failure ()
154   end
155 end
```

a successful CAS in line 3 and fails after completing the HELPIINSERT routine. In this case, the INSERT operation took effect, that is, the new key appears as a leaf of the tree, and a FIND(k) operation will return it. However, even though the operation was already linearized at the time of the crash, upon recovery process p is unaware of it. Moreover, looking for the new leaf in the tree is not helpful, since it might be k has been removed from the tree after the crash.

Moreover, if no recover is supplied, it may result an execution which is not well-formed. Consider for example the following scenario. A process p invoke an $Op_1 = \text{INSERT}(k_1)$ operation. After a successful CAS at line 3 the process crash. After recovering, p invoke an $Op_2 = \text{INSERT}(k_2)$ operation. Assume k_1 and k_2 belongs to a completely different parts of the tree. Then, p can complete inserting k_2 without having any affect on k_1 . Now, a process q performs FIND(k_1) which returns \perp , as the insertion of k_1 is not completed, follows by an FIND(k_2), which returns the leaf of k_2 . The INSERT(k_1) operation will be completed later by any process accessing 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. In the Linked-List algorithm it was enough to consider a marked node as if it has been already deleted. Nonetheless, the more complex BST implementation is even 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 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 complete the delete operation, and only then will return, according to whether the attempt was successful. The modified FIND routine is given in figure 6.

solved by having the FIND routine helping in case of a flagged or marked node. In such case, once a process flagged a node, then any operation, including FIND, will try and complete this operation, and thus either the operation is linearized at the time of the flagging, or there is no linearization at all. A more fine modification will be as follows. The routine FIND calls SEARCH and gets $\langle gp, p, l, pupdate, gpupdate \rangle$ in return. Then the process look to see if there is an unavoidable operation in progress that may affects its response. There are several cases to consider:

$l \rightarrow key \neq k$, then return FALSE if and only if $pupdate$ does not contains indication for inserting a key k .

$l \rightarrow key = k$, then if nor $gpupdate$ neither $pupdate$ contains indication for deleting l it is safe to return TRUE; Else, if $pupdate.state = \text{MARK}$ and $pupdate.info$'s Leaf is l , then return FALSE; Else, $gpupdate.state = \text{DFlag}$, and $gpupdate.info$'s Leaf is l . In such case, we do not know if the deletion will be completed or not (resulting a response FALSE or TRUE, accordingly), so first need to try and complete it. Hence, a HELPDELETE routine with $gpupdate$ is invoked. Once it completes, the response is whether the deletion took affect or not, that is, let $par := gpupdate.info \leftarrow p$, then it returns the evaluation of the statement $par \leftarrow update.state = \text{MARK}$.

The above solution prevents unwanted behaviour so that

However, such a solution is not efficient for the FIND routine, as it needs to complete other operations.

The helping mechanism, in which a process performing INSERT or DELETE first helps operation which marked or flagged the relevant nodes, guarantees the BST implementation is consistent even

in the presence of crash-recovery. Any operation has a single linearization point, which is the point where the successful CAS changes the tree data structure, that is, the CAS at the CAS-CHILD routine. Therefore, even in the case of a crash,

If a process p crash in a middle of executing an operation Op , then either Op did not flagged or marked any node, thus it did not took affect on the data-structure, nor it will in the future. However, if Op already flagged or marked a node, then the helping mechanism guarantees that

Shared variables: Node* head

Define Info: struct { *type*: OperationType, *pred, curr, new*: Node* }

Code for process p:

```
Procedure Insert(int key)
  Data: Node* pred, curr
         Node node = new Node (key)
47 while true do
48   <pred, curr> = lookup(key)
49   if curr.key == key then
50     return false
51   else
52     node.next = unmarked curr
53     Backup[p] = new Info (Insert, pred, curr, new)
54     if pred.next.CAS (unmarked curr, unmarked node) then
55       return true
56     end
57   end
58 end
```

```
Procedure Delete(int key)
  Data: Node* pred, curr, succ
59 while true do
60   <pred, curr> = lookup(key)
61   if curr.key != key then
62     return false
63   else
64     succ = curr.next
65     Backup[p] = new Info (Delete, pred, curr, null)
66     if curr.next.CAS (unmarked succ, marked succ) then
67       curr.deleter.CAS (null, p)
68       pred.next.CAS (unmarked curr, unmarked succ)
69       return (curr.deleter == p)
70     end
71   end
72 end
```

Figure 2: Recoverable Non-Blocking Linked-List

```

1  type Update {           ▷ stored in one CAS word
2      {CLEAN, DFlag, IFlag, MARK} state
3      Flag *info
4  }
5  type Internal {         ▷ subtype of Node
6       $Key \cup \{\infty_1, \infty_2\}$  key
7      Update update
8      Node *left, *right
9  }
10 type Leaf {             ▷ subtype of Node
11      $Key \cup \{\infty_1, \infty_2\}$  key
12 }
13 type IInfo {            ▷ subtype of Flag
14     Internal *p, *newInternal
15     Leaf *l
16 }
17 type DInfo {            ▷ subtype of Flag
18     Internal *gp, *p
19     Leaf *l
20     Update pupdate
21 }
▷ Initialization:
22 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.

```

23 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
24   Internal * $gp$ , * $p$ 
25   Node * $l := \text{Root}$ 
26   Update  $gpupdate, pupdate$                                      ▷ Each stores a copy of an update field
27   while  $l$  points to an internal node {
28        $gp := p$                                                      ▷ Remember parent of  $p$ 
29        $p := l$                                                        ▷ Remember parent of  $l$ 
30        $gpupdate := pupdate$                                            ▷ Remember update field of  $gp$ 
31        $pupdate := p \rightarrow \text{update}$                                ▷ Remember update field of  $p$ 
32       if  $k < l \rightarrow \text{key}$  then  $l := p \rightarrow \text{left}$  else  $l := p \rightarrow \text{right}$  ▷ Move down to appropriate child
33   }
34   return  $\langle gp, p, l, pupdate, gpupdate \rangle$ 
35 }

36 FIND(Key  $k$ ) : Leaf* {
37   Leaf * $l$ 
38    $\langle -, -, l, -, - \rangle := \text{SEARCH}(k)$ 
39   if  $l \rightarrow \text{key} = k$  then return  $l$ 
40   else return  $\perp$ 
41 }

42 INSERT(Key  $k$ ) : boolean {
43   Internal * $p, *newInternal$ 
44   Leaf * $l, *newSibling$ 
45   Leaf * $new :=$  pointer to a new Leaf node whose key field is  $k$ 
46   Update  $pupdate, result$ 
47   IInfo * $op$ 
48   while TRUE {
49        $\langle -, p, l, pupdate, - \rangle := \text{SEARCH}(k)$ 
50       if  $l \rightarrow \text{key} = k$  then return FALSE                             ▷ Cannot insert duplicate key
51       if  $pupdate.state \neq \text{CLEAN}$  then HELP( $pupdate$ )                 ▷ Help the other operation
52       else {
53            $newSibling :=$  pointer to a new Leaf whose key is  $l \rightarrow \text{key}$ 
54            $newInternal :=$  pointer to a new Internal node with key field  $\max(k, l \rightarrow \text{key})$ ,
                    update field  $\langle \text{CLEAN}, \perp \rangle$ , and with two child fields equal to  $new$  and  $newSibling$ 
                    (the one with the smaller key is the left child)
55            $op :=$  pointer to a new IInfo record containing  $\langle p, l, newInternal \rangle$ 
56            $result := \text{CAS}(p \rightarrow \text{update}, pupdate, \langle \text{IFlag}, op \rangle)$    ▷ iflag CAS
57           if  $result = pupdate$  then {                                   ▷ The iflag CAS was successful
58               HELPIINSERT( $op$ )                                         ▷ Finish the insertion
59               return TRUE
60           }
61           else HELP( $result$ )                                           ▷ The iflag CAS failed; help the operation that caused failure
62       }
63   }
64 }

65 HELPIINSERT(IInfo * $op$ ) {
    ▷ Precondition:  $op$  points to an IInfo record (i.e., it is not  $\perp$ )
66   CAS-CHILD( $op \rightarrow p, op \rightarrow l, op \rightarrow newInternal$ )             ▷ ichild CAS
67   CAS( $op \rightarrow p \rightarrow \text{update}, \langle \text{IFlag}, op \rangle, \langle \text{CLEAN}, op \rangle$ )   ▷ iunflag CAS
68 }

```

Figure 4: Pseudocode for SEARCH, FIND and INSERT.

```

69 DELETE(Key  $k$ ) : boolean {
70     Internal * $gp$ , * $p$ 
71     Leaf * $l$ 
72     Update  $pupdate$ ,  $gpupdate$ ,  $result$ 
73     DInfo * $op$ 

74     while TRUE {
75          $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
76         if  $l \rightarrow key \neq k$  then return FALSE ▷ Key  $k$  is not in the tree
77         if  $gpupdate.state \neq \text{CLEAN}$  then HELP( $gpupdate$ )
78         else if  $pupdate.state \neq \text{CLEAN}$  then HELP( $pupdate$ )
79         else {
80              $op :=$  pointer to a new DInfo record containing  $\langle gp, p, l, pupdate \rangle$  ▷ Try to flag  $gp$ 
81              $result := \text{CAS}(gp \rightarrow update, gpupdate, \langle DFlag, op \rangle)$  ▷ dflag CAS
82             if  $result = gpupdate$  then { ▷ CAS successful
83                 if HELPDELETE( $op$ ) then return TRUE ▷ Either finish deletion or unflag
84             }
85             else HELP( $result$ ) ▷ The dflag CAS failed; help the operation that caused the failure
86         }
87     }
88 }

89 HELPDELETE(DInfo * $op$ ) : boolean {
90     ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
91     Update  $result$  ▷ Stores result of mark CAS
92      $result := \text{CAS}(op \rightarrow p \rightarrow update, op \rightarrow pupdate, \langle \text{MARK}, op \rangle)$  ▷ mark CAS
93     if  $result = op \rightarrow pupdate$  or  $result = \langle \text{MARK}, op \rangle$  then { ▷  $op \rightarrow p$  is successfully marked
94         HELPMARKED( $op$ ) ▷ Complete the deletion
95         return TRUE ▷ Tell DELETE routine it is done
96     }
97     else { ▷ The mark CAS failed
98         HELP( $result$ ) ▷ Help operation that caused failure
99          $\text{CAS}(op \rightarrow gp \rightarrow update, \langle DFlag, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ backtrack CAS
100         return FALSE ▷ Tell DELETE routine to try again
101     }

102 HELPMARKED(DInfo * $op$ ) {
103     ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
104     Node * $other$ 
105     ▷ Set  $other$  to point to the sibling of the node to which  $op \rightarrow l$  points
106     if  $op \rightarrow p \rightarrow right = op \rightarrow l$  then  $other := op \rightarrow p \rightarrow left$  else  $other := op \rightarrow p \rightarrow right$ 
107     ▷ Splice the node to which  $op \rightarrow p$  points out of the tree, replacing it by  $other$ 
108      $\text{CAS-CHILD}(op \rightarrow gp, op \rightarrow p, other)$  ▷ dchild CAS
109      $\text{CAS}(op \rightarrow gp \rightarrow update, \langle DFlag, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ dunflag CAS
110 }

111 HELP(Update  $u$ ) { ▷ General-purpose helping routine
112     ▷ Precondition:  $u$  has been stored in the  $update$  field of some internal node
113     if  $u.state = \text{IFlag}$  then HELPIINSERT( $u.info$ )
114     else if  $u.state = \text{MARK}$  then HELPMARKED( $u.info$ )
115     else if  $u.state = \text{DFlag}$  then HELPDELETE( $u.info$ )
116 }

117 CAS-CHILD(Internal * $parent$ , Node * $old$ , Node * $new$ ) {
118     ▷ Precondition:  $parent$  points to an Internal node and  $new$  points to a Node (i.e., neither is  $\perp$ )
119     ▷ This routine tries to change one of the child fields of the node that  $parent$  points to from  $old$  to  $new$ .
120     if  $new \rightarrow key < parent \rightarrow key$  then
121          $\text{CAS}(parent \rightarrow left, old, new)$ 
122     else
123          $\text{CAS}(parent \rightarrow right, old, new)$ 
124 }

```

Figure 5: Pseudocode for DELETE and some auxiliary routines.

```

FIND(Key  $k$ ) : Leaf* {
119   Internal *gp, *p
120   Leaf *l
121   Update pupdate, gpupdate

122    $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
123   if  $l \rightarrow key \neq k$  then {
124       if ( $pupdate.state = \text{IFlag}$  and  $pupdate.info$  attempt to add key  $k$ ) then return  $pupdate.info \leftarrow key\ k$ 
125       else return  $\perp$ 
126   }
127   if ( $pupdate.state = \text{MARK}$  and  $pupdate.info \leftarrow l \leftarrow key = k$ ) then return  $\perp$ 
128   if ( $gpupdate.state = \text{DFlag}$  and  $gpupdate.info \leftarrow l \leftarrow key = k$ ) then {
129        $op := gpupdate.info$ 
130        $result := \text{CAS}(op \rightarrow p \rightarrow update, op \rightarrow pupdate, \langle \text{MARK}, op \rangle)$   $\triangleright$  mark CAS
131       if ( $result = op \rightarrow pupdate$  or  $result = \langle \text{MARK}, op \rangle$ ) then return  $\perp$   $\triangleright$   $op \rightarrow p$  is successfully marked
132   }
133   return  $l$ 
134 }

```

Figure 6: R-linearizability FIND routine

```

1  type Update {  $\triangleright$  stored in one CAS word
2      {CLEAN, DFlag, IFlag, MARK} state
3      Flag *info
4  }
5  type Internal {  $\triangleright$  subtype of Node
6       $Key \cup \{\infty_1, \infty_2\}$  key
7      Update update
8      Node *left, *right
9  }
10 type Leaf {  $\triangleright$  subtype of Node
11      $Key \cup \{\infty_1, \infty_2\}$  key
12 }
13 type IInfo {  $\triangleright$  subtype of Flag
14     Internal *p, *newInternal
15     Leaf *l
16     boolean complete
17 }
18 type DInfo {  $\triangleright$  subtype of Flag
19     Internal *gp, *p
20     Leaf *l
21     Update pupdate
22     boolean complete
23 }
24  $\triangleright$  Initialization:
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 7: Type definitions and initialization.

```

25 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  $gpupdate$ 
26 Internal * $gp$ , * $p$ 
27 Node * $l := \text{Root}$ 
28 Update  $gpupdate, pupdate$  ▷ Each stores a copy of an update field
29 while  $l$  points to an internal node {
30      $gp := p$  ▷ Remember parent of  $p$ 
31      $p := l$  ▷ Remember parent of  $l$ 
32      $gpupdate := pupdate$  ▷ Remember update field of  $gp$ 
33      $pupdate := p \rightarrow \text{update}$  ▷ Remember update field of  $p$ 
34     if  $k < l \rightarrow \text{key}$  then  $l := p \rightarrow \text{left}$  else  $l := p \rightarrow \text{right}$  ▷ Move down to appropriate child
35 }
36 return  $\langle gp, p, l, pupdate, gpupdate \rangle$ 
37 }

38 FIND(Key  $k$ ) : Leaf* {
39     Leaf * $l$ 
40      $\langle -, -, l, -, - \rangle := \text{SEARCH}(k)$ 
41     if  $l \rightarrow \text{key} = k$  then return  $l$ 
42     else return  $\perp$ 
43 }

44 INSERT(Key  $k$ ) : boolean {
45     Internal * $p$ , * $\text{newInternal}$ 
46     Leaf * $l$ , * $\text{newSibling}$ 
47     Leaf * $\text{new} :=$  pointer to a new Leaf node whose key field is  $k$ 
48     Update  $pupdate, \text{result}$ 
49     IInfo * $op$ 
50     while TRUE {
51          $\langle -, p, l, pupdate, - \rangle := \text{SEARCH}(k)$ 
52         if  $l \rightarrow \text{key} = k$  then return FALSE ▷ Cannot insert duplicate key
53         if  $pupdate.state \neq \text{CLEAN}$  then HELP( $pupdate$ ) ▷ Help the other operation
54         else {
55              $\text{newSibling} :=$  pointer to a new Leaf whose key is  $l \rightarrow \text{key}$ 
56              $\text{newInternal} :=$  pointer to a new Internal node with key field  $\max(k, l \rightarrow \text{key})$ ,
                    update field  $\langle \text{CLEAN}, \perp \rangle$ , and with two child fields equal to  $\text{new}$  and  $\text{newSibling}$ 
                    (the one with the smaller key is the left child)
57              $op :=$  pointer to a new IInfo record containing  $\langle p, l, \text{newInternal}, \text{FALSE} \rangle$ 
58              $\text{Announce}[p] := op$ 
59              $\text{result} := \text{CAS}(p \rightarrow \text{update}, pupdate, \langle \text{IFlag}, op \rangle)$  ▷ iflag CAS
60             if  $\text{result} = pupdate$  or  $\text{result} = \langle \text{IFlag}, op \rangle$  then { ▷ The iflag CAS was successful
61                 HELPINSERT( $op$ ) ▷ Finish the insertion
62                 return TRUE
63             }
64             else HELP( $\text{result}$ ) ▷ The iflag CAS failed; help the operation that caused failure
65         }
66         if  $op \rightarrow \text{complete} = \text{TRUE}$  then
67             return TRUE
68     }
69 }

70 HELPINSERT(IInfo * $op$ ) {
    ▷ Precondition:  $op$  points to an IInfo record (i.e., it is not  $\perp$ )
71     CAS-CHILD( $op \rightarrow p, op \rightarrow l, op \rightarrow \text{newInternal}$ ) ▷ ichild CAS
72     CAS( $op \rightarrow p \rightarrow \text{update}, \langle \text{IFlag}, op \rangle, \langle \text{CLEAN}, op \rangle$ ) ▷ iunflag CAS
73      $op \rightarrow \text{complete} := \text{TRUE}$  ▷ mark the operation as completed
74 }

```

Figure 8: Pseudocode for SEARCH, FIND and INSERT.

```

75 DELETE(Key  $k$ ) : boolean {
76   Internal * $gp$ , * $p$ 
77   Leaf * $l$ 
78   Update  $pupdate$ ,  $gpupdate$ ,  $result$ 
79   DInfo * $op$ 
80   while TRUE {
81      $\langle gp, p, l, pupdate, gpupdate \rangle := \text{SEARCH}(k)$ 
82     if  $l \rightarrow key \neq k$  then return FALSE ▷ Key  $k$  is not in the tree
83     if  $gpupdate.state \neq \text{CLEAN}$  then HELP( $gpupdate$ )
84     else if  $pupdate.state \neq \text{CLEAN}$  then HELP( $pupdate$ )
85     else { ▷ Try to flag  $gp$ 
86        $op :=$  pointer to a new DInfo record containing  $\langle gp, p, l, pupdate, \text{FALSE} \rangle$ 
87        $\text{Announce}[p] := op$ 
88        $result := \text{CAS}(gp \rightarrow update, gpupdate, \langle \text{DFlag}, op \rangle)$  ▷ dflag CAS
89       if  $result = gpupdate$  or  $result = \langle \text{DFlag}, op \rangle$  then { ▷ CAS successful
90         if HELPDELETE( $op$ ) then return TRUE ▷ Either finish deletion or unflag
91       }
92       else HELP( $result$ ) ▷ The dflag CAS failed; help the operation that caused the failure
93     }
94     if  $op \rightarrow complete = \text{TRUE}$  then
95       return TRUE
96   }
97 }

98 HELPDELETE(DInfo * $op$ ) : boolean {
99   ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
100   Update  $result$  ▷ Stores result of mark CAS
101    $result := \text{CAS}(op \rightarrow p \rightarrow update, op \rightarrow pupdate, \langle \text{MARK}, op \rangle)$  ▷ mark CAS
102   if  $result = op \rightarrow pupdate$  or  $result = \langle \text{MARK}, op \rangle$  then { ▷  $op \rightarrow p$  is successfully marked
103     HELPMARKED( $op$ ) ▷ Complete the deletion
104     return TRUE ▷ Tell DELETE routine it is done
105   }
106   else { ▷ The mark CAS failed
107     HELP( $result$ ) ▷ Help operation that caused failure
108      $\text{CAS}(op \rightarrow gp \rightarrow update, \langle \text{DFlag}, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ backtrack CAS
109     return FALSE ▷ Tell DELETE routine to try again
110   }

111 HELPMARKED(DInfo * $op$ ) {
112   ▷ Precondition:  $op$  points to a DInfo record (i.e., it is not  $\perp$ )
113   Node * $other$ 
114   ▷ Set  $other$  to point to the sibling of the node to which  $op \rightarrow l$  points
115   if  $op \rightarrow p \rightarrow right = op \rightarrow l$  then  $other := op \rightarrow p \rightarrow left$  else  $other := op \rightarrow p \rightarrow right$ 
116   ▷ Splice the node to which  $op \rightarrow p$  points out of the tree, replacing it by  $other$ 
117    $\text{CAS-CHILD}(op \rightarrow gp, op \rightarrow p, other)$  ▷ dchild CAS
118    $\text{CAS}(op \rightarrow gp \rightarrow update, \langle \text{DFlag}, op \rangle, \langle \text{CLEAN}, op \rangle)$  ▷ dunflag CAS
119    $op \rightarrow complete := \text{TRUE}$  ▷ mark the operation as completed
120 }

121 HELP(Update  $u$ ) { ▷ General-purpose helping routine
122   ▷ Precondition:  $u$  has been stored in the  $update$  field of some internal node
123   if  $u.state = \text{IFlag}$  then HELPINSERT( $u.info$ )
124   else if  $u.state = \text{MARK}$  then HELPMARKED( $u.info$ )
125   else if  $u.state = \text{DFlag}$  then HELPDELETE( $u.info$ )
126 }

127 CAS-CHILD(Internal * $parent$ , Node * $old$ , Node * $new$ ) {
128   ▷ Precondition:  $parent$  points to an Internal node and  $new$  points to a Node (i.e., neither is  $\perp$ )
129   ▷ This routine tries to change one of the child fields of the node that  $parent$  points to from  $old$  to  $new$ .
130   if  $new \rightarrow key < parent \rightarrow key$  then
131      $\text{CAS}(parent \rightarrow left, old, new)$ 
132   else
133      $\text{CAS}(parent \rightarrow right, old, new)$ 
134 }

```

Figure 9: Pseudocode for DELETE and some auxiliary routines.