# 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 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 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 = $ INSERT($k_1$) operation. $q$ performs a successful CAS in line 96 followed by a crush. After recovering, $q$ invoke an $Op_2 = $ 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 FIND($k_1$) which returns $\bot$, as the insertion of $k_1$ is not completed, followed by FIND($k_2$), which returns the leaf of $k_2$. The 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 drown 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 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-linearzability. 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. We assume for simplicity nodes and Flag records are always allocated new memory locations, although it is enough to require no location is reallocated as long as there is a chain of pointers leading to it. 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 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 $q$ from some execution of the original algorithm in which more processes can participate, then it has to return the same response on both. The proof relies on several key arguments given below.

1. 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

any changes to the BST, but rather 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.

2. If an internal node $nd_1$ stops pointing to a node $nd_2$ at some point of the execution, it can not point to $nd_2$ again. This attributes to the fact an INSERT presents a node with two new children. Therefore, if $nd_2$ is a leaf, it can either be deleted, or replace with a new copy when an insert operation takes affect. Otherwise, $nd_2$ is an internal node, and as such, it can not be replaced by an insert operation (which only allows to replace a leaf), and therefore it can only be removed from the tree.

3. The field update of a node $nd$ can have any value only once along the execution. This follows from the fact that any attempt to perform an operation creates new record in the memory. If $nd \leftarrow update$ is marked, it can be unmarked or change. Otherwise, any attempt to flag it uses a new created record $op$. If the attempt succeed, then eventually it will be unflagged while still referring to $op$. In order to replace the value again, there must be an operation reading $nd \leftarrow update$ after it was unflagged (as any operation first help a flagged node). This operation must create a new record, and thus we can use the same argument. As a corollary, if a process successfully flag or mark a node, there was no change in the node since the last time it read the node update field.

Assume process $q$ performs an INSERT(k) operation. As argued in argument 1, a crash before writing to Announce[q] implies no changes has been made to the BST. Assume thus $q$ executes line 95, i.e., $q$ stored in 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 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 Announce[q] without having any affect on the BST. Hence, the same argument can be applied once $q$ writes to Announce[q] again.

Assume $q$ performs a successful CAS in line 96. Then, a reference to the IInfo $op$ is stored in $p \leftarrow update$, which is also flagged. Following argument 3, $p$ was not changed since the SEARCH routine read it, and it still points to $l$. Starting from this point, no changes can be made to $p$, except for the change point to by $op$, as the node is flagged. Now, only the first process

Relying on the correctness of the original algorithm, no matter how many times HELPINSERT(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 := $ SEARCH($k$)
6            if *gpupdate.state* $\neq$ CLEAN then HELP(*gpupdate*)
7            else if *pupdate.state* $\neq$ CLEAN then 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 := $ 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($op \rightarrow p \rightarrow update, op \rightarrow pupdate, \langle$MARK$, op\rangle$)     ▷ **mark** CAS
27            if ($result = op \rightarrow pupdate$ or $result = \langle$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 ∪ {∞₁, ∞₂} key
37       Update update
38       Node *left, *right
39   }
40   type Leaf {                ▷ subtype of Node
41       Key ∪ {∞₁, ∞₂} key
42   }
43   type IInfo {               ▷ subtype of Flag
44       Internal *p, *newInternal
45       Leaf *l
46       {CLEAN,DFlag,IFlag,MARK, TRUE} done
47   }
48   type DInfo {               ▷ subtype of Flag
49       Internal *gp, *p
50       Leaf *l
51       Update pupdate
52       {CLEAN,DFlag,IFlag,MARK, TRUE} done
53   }
     ▷ Initialization:
54   shared Internal *Root := pointer to new Internal node
         with key field ∞₂, update field ⟨CLEAN, ⊥⟩, and
         pointers to new Leaf nodes with keys ∞₁ and
         ∞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 126
61   }
```

Figure 4: RECOVER routine

62  SEARCH($Key\ k$) : $\langle$Internal\*, Internal\*, Leaf\*, Update, Update$\rangle$ {
        ▷ 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 $p \to left$ has contained $l$ (if $k < p \to key$) or $p \to right$ has contained $l$ (if $k \geq p \to key$)
        ▷ (3) $p \to update$ has contained *pupdate*
        ▷ (4) if $l \to key \neq \infty_1$, then the following three statements hold:
        ▷     (4a) *gp* points to an Internal node
        ▷     (4b) either $gp \to left$ has contained $p$ (if $k < gp \to key$) or $gp \to right$ has contained $p$ (if $k \geq gp \to key$)
        ▷     (4c) $gp \to update$ has contained *gpupdate*
63      Internal \*$gp$, \*$p$
64      Node \*$l$ := $Root$
65      Update *gpupdate*, *pupdate*                                    ▷ Each stores a copy of an *update* field

66      while $l$ points to an internal node {
67          $gp$ := $p$                                                 ▷ Remember parent of $p$
68          $p$ := $l$                                                  ▷ Remember parent of $l$
69          *gpupdate* := *pupdate*                                     ▷ Remember *update* field of $gp$
70          *pupdate* := $p \to update$                                 ▷ Remember *update* field of $p$
71          if $k < l \to key$ then $l := p \to left$ else $l := p \to right$   ▷ Move down to appropriate child
72      }
73      return $\langle gp, p, l, pupdate, gpupdate\rangle$
74  }

75  FIND($Key\ k$) : Leaf\* {
76      Leaf \*$l$

77      $\langle -, -, l, -, -\rangle$ := SEARCH($k$)
78      if $l \to key = k$ then return $l$
79      else return $\bot$
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          $\langle -, p, l, pupdate, -\rangle$ := SEARCH($k$)
89          if $l \to key = k$ then return FALSE                        ▷ Cannot insert duplicate key
90          if $pupdate.state \neq$ CLEAN then HELP($pupdate$)           ▷ Help the other operation
91          else {
92              $newSibling$ := pointer to a new Leaf whose key is $l \to key$
93              $newInternal$ := pointer to a new Internal node with *key* field $\max(k, l \to key)$,
                        *update* field $\langle$CLEAN, $\bot\rangle$, and with two child fields equal to $new$ and $newSibling$
                        (the one with the smaller key is the left child)
94              $op$ := pointer to a new IInfo record containing $\langle p, l, newInternal,$ CLEAN$\rangle$
95              $Announce[id]$ := $op$
96              $result$ := CAS($p \to update, pupdate, \langle$IFlag, $op\rangle$)   ▷ **iflag** CAS
97              if $result = pupdate$ or $result = \langle$IFlag, $op\rangle$ then {   ▷ The iflag CAS was successful
98                  HELPINSERT($op$)                                    ▷ Finish the insertion
99                  return TRUE
100             }
101             else HELP($result$)           ▷ The iflag CAS failed; help the operation that caused failure
102         }
103         if $op \to done =$ TRUE then
104             return TRUE
105     }
106 }

107 HELPINSERT(IInfo \*$op$) {
        ▷ *Precondition*: $op$ points to an IInfo record (*i.e.*, it is not $\bot$)
108     $op \to done$ := IFlag                                          ▷ announce the iflag CAS was successful
109     CAS-CHILD($op \to p, op \to l, op \to newInternal$)             ▷ **ichild** CAS
110     CAS($op \to p \to update, \langle$IFlag, $op\rangle, \langle$CLEAN, $op\rangle$)   ▷ **iunflag** CAS
111     $op \to done$ := TRUE                                           ▷ announce the operation is completed
112 }

Figure 5: Pseudocode for SEARCH, FIND and INSERT.

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

118      while TRUE {
119          ⟨gp, p, l, pupdate, gpupdate⟩ := SEARCH(k)
120          if l → key ≠ k then return FALSE                                    ▷ Key k is not in the tree
121          if gpupdate.state ≠ CLEAN then HELP(gpupdate)
122          else if pupdate.state ≠ CLEAN then HELP(pupdate)
123          else {                                                              ▷ Try to flag gp
124              op := pointer to a new DInfo record containing ⟨gp, p, l, pupdate, CLEAN⟩
125              Announce[id] := op
126              result := CAS(gp → update, gpupdate, ⟨DFlag, op⟩)               ▷ dflag CAS
127              if result = gpupdate or result = ⟨DFlag, op⟩ then {             ▷ CAS successful
128                  if HELPDELETE(op) then return TRUE                         ▷ Either finish deletion or unflag
129              }
130              else HELP(result)              ▷ The dflag CAS failed; help the operation that caused the failure
131          }
132          if op → done = TRUE then
133              return TRUE
134      }
135  }

136  HELPDELETE(DInfo *op) : boolean {
         ▷ Precondition: op points to a DInfo record (i.e., it is not ⊥)
137      Update result                                                          ▷ Stores result of mark CAS

138      result := CAS(op → p → update, op → pupdate, ⟨MARK, op⟩)               ▷ mark CAS
139      if result = op → pupdate or result = ⟨MARK, op⟩ then {                 ▷ op → 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          CAS(op → gp → update, ⟨DFlag, op⟩, ⟨CLEAN, op⟩)                    ▷ backtrack CAS
146          return FALSE                                                       ▷ Tell DELETE routine to try again
147      }
148  }

149  HELPMARKED(DInfo *op) {
         ▷ Precondition: op points to a DInfo record (i.e., it is not ⊥)
150      Node *other

         ▷ Set other to point to the sibling of the node to which op → l points
151      if op → p → right = op → l then other := op → p → left else other := op → p → right
         ▷ Splice the node to which op → p points out of the tree, replacing it by other
152      CAS-CHILD(op → gp, op → p, other)                                      ▷ dchild CAS
153      CAS(op → gp → update, ⟨DFlag, op⟩, ⟨CLEAN, op⟩)                        ▷ dunflag CAS
154      op → done := TRUE                                                      ▷ mark the operation as completed
155  }

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

161  CAS-CHILD(Internal *parent, Node *old, Node *new) {
         ▷ Precondition: parent points to an Internal node and new points to a Node (i.e., neither is ⊥)
         ▷ This routine tries to change one of the child fields of the node that parent points to from old to new.
162      if new → key < parent → key then
163          CAS(parent → left, old, new)
164      else
165          CAS(parent → right, old, new)
166  }
```

Figure 6: Pseudocode for DELETE and some auxiliary routines.