

1 A relaxed balanced lock-free binary search tree

2 Manish Singh¹, Lindsay Groves², and Alex Potanin²

3 ¹ Wellington Institute of Technology

4 ² Victoria University of Wellington

5 manish.singh@weltec.ac.nz

6 {lindsay.groves, alex.potanin}@ecs.vuw.ac.nz

7 **Abstract.** This paper presents a new relaxed balanced concurrent bi-
8 nary search tree using a single word compare and swap primitive, in
9 which all operations are lock-free. Our design separates balancing ac-
10 tions from update operations and includes a lock-free balancing mech-
11 anism in addition to the insert, search, and relaxed delete operations.
12 Search in our design is not affected by ongoing concurrent update op-
13 erations or by the movement of nodes by tree restructuring operations.
14 Our experiments show that our algorithm performs better than other
15 state-of-the-art concurrent BSTs.

16 **Keywords:** Concurrent Data Structure · Lock-Free · Binary Search Tree.

17 1 Introduction

18 Recently, the chip manufacturing company, AMD, has released the Threadripper
19 CPU with 64 cores and 128 threads for desktops. As CPU manufacturers embed
20 more and more cores in their multi-core processors, the need to design scalable
21 and efficient concurrent data structures has also increased. Considerable research
22 has been done towards making both blocking and non-blocking concurrent ver-
23 sions of sequential data structures. Unlike blocking a concurrent data structure
24 design is non-blocking, or lock-free, if it ensures at any time at least one thread
25 can complete its operation. Performance has always been an important factor
26 and will drive the design of these data structures.

27 The Binary Search Tree (BST) is a fundamental data structure. Many concu-
28 rent BST, both blocking and non-blocking, have been proposed [1, 4–6, 14, 15, 17].
29 However, only a few designs include self-balancing operations. Most of the pub-
30 lished work tries to emulate a sequential implementation. This results in per-
31 formance compromise as strict sequential invariants must be maintained. In a
32 concurrent environment the effect of some operations might cancel out the ef-
33 fects of other operations. In the case of a self-balancing AVL tree each insert or
34 delete operation requires a balancing operation to be performed immediately to
35 preserve the height-balanced property. A balancing operation might cancel out
36 the effects of other balancing operations on the same set of nodes.

37 We have designed a lock-free relaxed balanced BST, hereafter referred as
38 RBLFBST, in which balancing operations are decoupled from the regular in-

sert and delete operations¹. The idea of relaxing some structural properties of
 a sequential data structure in their concurrent version has been tried in sev-
 eral previous works. Lazy deletion [9, 10] is an example of relaxing the require-
 ment that nodes are immediately removed from the data structure. In a relaxed
 self-balancing such as AVL [11], chromatic trees [16], rotation operations, are
 performed separately from the insertion and deletion operation.

2 Related Work

The first lock-free BST dates back to the 90’s [19] where the author suggested a
 design based on a threaded binary tree representation, but did not discuss the
 implementation of his design. Fraser et al [7] described a lock-free implementa-
 tion of an internal BST using multiple CAS (MCAS) to update multiple parts
 of the tree structure atomically. This algorithm uses a threaded representation
 of a BST. A delete operation required up to 8 memory locations to be updated
 atomically which added an appreciable overhead to the design. Based on the
 cooperative technique of [2], Ellen et al. [6] developed the first specific lock-free
 external BST algorithm, where internal nodes only act as routing nodes, which
 simplifies the deletion operation.

Using a helping technique similar to [6], Howley and Jones [17] presented a
 non-blocking algorithm for internal BSTs. They added an extra field *operation*
 to each node of the BST, in which all the information related to a particular up-
 date operation can be stored. A delete operation in this work could use as many
 as 9 CAS instructions. However, the paper gives experimental results in which
 this design outperforms [6, 15]. Aravind et al. [14] presented a lock-free external
 BST similar to [6], but using edge-marking. Being a leaf-oriented tree this design
 also has a smaller contention window for insert and delete operations and uses
 fewer auxiliary nodes to coordinate among conflicting operations. Another con-
 temporary paper [18] describes a general template for non-blocking trees. This
 template uses multi-word versions of LL/SC and VL primitives, making it easier
 to update multiple parts of the tree. The paper also presented a fine-grained
 synchronized version of a leaf-oriented chromatic tree based on their template.

The concurrent AVL tree of Bronson et al. [15] uses lock coupling which uses
 per-node locks to remove conflicts between concurrent updates, and a relaxed
 balancing property which does not enforce the strict AVL property. This design
 uses a partially external tree such that by default nodes behave as per an internal
 tree but in order to simplify the removal of nodes with two children, the removed
 node objects are allowed to remain in the structure and act as routing nodes.
 This avoids the problem of locking large parts of the tree if a delete operation was
 implemented exactly as in a sequential BST. These deleted (logically) nodes can
 then be either removed during a later re-balancing operation or, if the node’s key
 is reinserted, the key can be made part of the set again. This partially external
 tree design was experimentally shown to have a small increase (0-20%) on the

¹ A poster describing the design of RBLFBST was presented in ICPP 2019, Kyoto, Japan

80 total number of nodes required to contain the key set. Crain et al. [4] present a
 81 lock-based tree in which balancing operations are decoupled from the insert and
 82 delete, and are done by a separate dedicated thread. Keys that are deleted are
 83 not immediately removed from the tree but are only marked as deleted as in [15].
 84 Later, a dedicated thread can remove nodes with deleted keys that have single
 85 or no child. The balancing mechanism used closely mirrors that in [3]. Despite
 86 claiming performance improvement by more than double to that of [15] another
 87 lock-based design this tree still is lock-based. More recently, Drachler et al. [5]
 88 proposed a lock-based internal BST supporting wait-free search operations and
 89 lock-based update operations via logical ordering. Their design uses a similar
 90 threaded binary tree representation as in [19].

91 All the designs of concurrent BST with balancing operations described above
 92 use locks to synchronize concurrent updates and therefore are not immune to
 93 problems that are associated with locking in general. Based on techniques used in
 94 previous research we present a concurrent BST in which all the update operations
 95 are lock-free. To our knowledge, our design is the first AVL tree based lock-free
 96 partially external BST which includes balancing operation for all the update
 97 operations.

98 3 Algorithm Description

99 3.1 Overview

100 We implement a set using a BST. Our implementation supports concurrent ex-
 101 ecutions of $search(k)$: to determine whether the key k is in the set, $insert(k)$: to
 102 add the key k to the set, $delete(k)$: to remove the key k from the set. To ensure
 103 that the tree does not become unbalanced, causing the operations to have linear
 104 cost, our design supports a relaxed tree balancing mechanism.

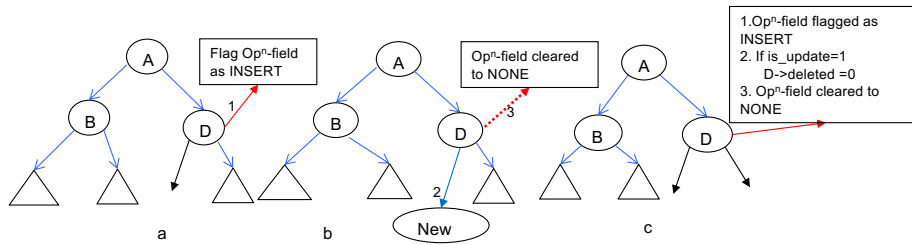


Fig. 1: left child pointer of Node D is identified as an insertion point for a 'new' key. (a). An *operation-descriptor* flagged as INSERT is inserted in to the *operation* field of node D (1). (b). After successful insertion of *operation-descriptor*, the node having key 'new' as a child of node D is physically inserted tiny(2) and then the *operation* field of node D is cleared by setting its flag to NONE (3). (c) Shows the case when the key to be inserted is found in the tree but is deleted. In this case, the *deleted* field is set to false.

105 Any thread attempting to insert a *key*, upon finding the insertion point first
 106 announces its intention. The thread first collects the information required in
 107 an *operation-descriptor* and tries to insert it in the *operation* field of the node
 108 identified as the insertion point (the parent of the new node). Once the *operation-*
 109 *descriptor* is inserted into the node the *key* is considered to be part of the set.

110 Then the new node is inserted using a CAS to the appropriate child pointer of
 111 the parent node as shown in figure 1(a) and (b).

112 To remove a key, first the node having the key to be removed is flagged as
 113 deleted (*deleted* bit is set to true) and then physically removed later. Once the
 114 *deleted* bit is set the key is considered to be removed (logically) from the tree.
 115 The physical removal is started by a separate maintenance thread by marking
 116 the deleted nodes having at most one child. Once marked any other thread can
 117 physically remove the node from the tree. Deleted nodes with two children are
 118 not physically removed from the tree until they have less than two children. This
 119 relaxed deletion is done to reduce contention in the tree as the delete algorithm
 120 does not need to locate and update the node to be deleted and the replacement
 121 node, which might involve several restarts. This relaxed deletion approach also
 122 means that a thread doing an insert operation may find the intended key already
 123 in the tree but the node is flagged as deleted. In this case, the key can be made
 124 part of the set again by setting the *deleted* bit off.

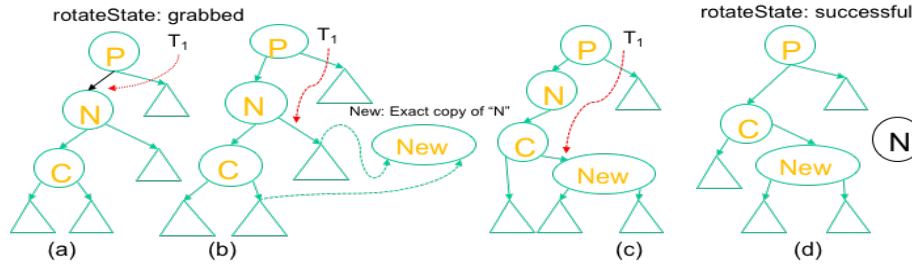


Fig. 2: Right rotation example (similarly for the left). (a) Shows the case where the maintenance thread has found a balance condition violation at node N and has successfully inserted a rotation *operation-descriptor* in the *operation* field of the parent of N, P, N, and child of N, C. After that a *new* node having key and other fields exactly as the node N is created. (b) Then the right and the left pointer of the *new node* are allocated to those children which the node N would have got after the rotation. (c) The next step is to insert the *new* node to its position after *rotation*. In this case, the *new* node would be the right-child of node C. (d) The third and last step is to connect node the parent P to child C effectively removing node N from the tree. A thread T1 carrying out search is oblivious to the movement of nodes by rotations

125 Our balancing mechanism is based on heights of nodes and rotation opera-
 126 tions, as in sequential AVL trees. However, balancing is carried out separately to
 127 the update operations, similar to that of [4, 12]. The balancing adjustments are
 128 performed by the maintenance thread. If the balancing condition at any node
 129 is found to be violated then the maintenance thread collects all the information
 130 required to do rotation operation in to an *operation-descriptor* and tries to insert
 131 it in to the *operation* field of the parent of the node. After successfully inserting
 132 the operation descriptor, the maintenance thread (or any other thread) tries to
 133 do the same for the node and the appropriate child of the node involved in the
 134 rotation. Figure 2 shows the remaining steps for the rotation operation (right ro-
 135 tation) once the *operation-descriptor* is inserted in all the three nodes involved.
 136 Due to the ongoing concurrent updates, it is difficult to determine the exact
 137 height of a node at any point in the execution. Thus, the balancing condition
 138 in RBLFBST is based on apparent local heights of the node and its children. It

139 should also be noted that the rotation technique shown in figure 2 makes ongoing
 140 concurrent search operations oblivious to any node moving up or down the tree.

```

1  struct node_t {
2      int key
3      node_t* left, right
4      operation_t* op
5      int local_height, lh, rh
6      bool deleted, removed
7  }
8  union operation_t {
9      insert_op_t insert_op
10     rotate_op_t rotate_op
11 }
12 struct insert_op_t {
13     bool_t is_left
14     bool_t is_update=FALSE
15     node_t* expected
16     node_t* new
17 } insert_op_t
18 struct rotate_op_t {
19     int /*volatile*/ state =
20         UNDECIDED
21     node_t* parent
22     node_t* node
23     node_t* child
24     operation_t* pOp
25     operation_t* nOp
26     operation_t* cOp
27     bool rightR
28     bool dir
29 } rotate_op_t
29 //operation marking status
30 NONE 0, MARK 1, ROTATE 2,
    INSERT 3

```

Algorithm 1: Basic structures used

```

Input: int k, node_t* root
Output: true, if the key is found in
           the set otherwise false
1  node_t* node, next
2  operation_t* node_op
3  bool result = FALSE
4  node = root
5  node_op = node→op
6  int node_key
7  next = node→right
8  while !ISNULL(next) do
9      node = next
10     node_op = node→op
11     node_key = node→key
12     if k < node_key then
13         next = node→left
14     else if k > node_key then
15         next = node→right
16     else
17         result = TRUE
18         break
19 if result && node→deleted then
20     if GETFLAG(node→op) ==
21         INSERT then
22         if
23             node_op→insert_op.new→key
24             == k then
25             return TRUE
26         return FALSE
27 return result

```

Algorithm 2: Search

143 All the operations in RBLFBST are lock-free and linearisable. Once an operation
 144 is flagged with the corresponding FLAG and the *operation-descriptor* is
 145 inserted, it is considered to be logically completed. Any other thread can complete
 146 the announced operation using the information available in the *operation-*
 147 *descriptor*. Updating multiple locations using a single word CAS while preserving
 148 atomicity is a challenging task. Especially, in case of a rotation operation, three
 149 locations are needed to be updated atomically. We achieve this by careful design
 150 of our operation, described in detail in later sections. Another notable feature
 151 of this design is that we managed to keep the most frequently used operation in
 152 BSTs, *search*, free of additional coordination.

153 3.2 Detailed Algorithm

154 **Structures:** Algorithm 1 shows various structures used in the implementation.
 155 As in a sequential BST, a node has key, right, and left pointers to its correspond-
 156 ing child which are set to NULL for every new node. To synchronize various
 157 concurrent operations, each node has an operation pointer field, *op*, that is used
 158 to announce the intended operation and contains a pointer to that particular
 159 *operation-descriptor*. A node in RBLFBST has 3 heights, *local-height*: updated

160 by adding one to the maximum *local-height* of its children, *right-height*: *local-*
 161 *height* of right child, and *left-height*: *local-height* of left child, which are stored in
 162 fields *local-heights*, *lh* and *rh* respectively. In addition to that it has two boolean
 163 fields *deleted* and *removed* which are initialized to false. If a node has its *deleted*
 164 field set that means it is deleted from the set but could still be in the tree, while
 165 the *removed* field is set when the node is removed from the tree.

<pre> int seek (int key, node_t** parent, operation_t** parent_op, node_t** node, operation_t** node_op, node_t* aux_root, node_t* root) 1 int result, nodekey 2 node_t* next 3 retry: 4 result = NOT_FOUND_L 5 *node = aux_root 6 *node_op = *node→op 7 if GETFLAG(*node_op) ≠ NONE then 8 if auxroot == root then 9 bst_help_insert((operation_t*) UNFLAG(*node_op), *node) 10 goto retry 11 next = (nodelt*) (*node)→right 12 while !ISNULL(next) do 13 *parent = *node 14 *parent_op = *node_op 15 *node = next 16 *node_op = (*node)→op 17 nodekey = (*node)→key 18 if key < nodekey then 19 result = NOT_FOUND_L 20 next = (*node)→left 21 else if key > nodekey then 22 result = NOT_FOUND_R 23 next = (*node)→right 24 else 25 result = FOUND 26 break 27 if GETFLAG(*node_op)≠ NONE then 28 help(*parent, *parent_op, *node, *node_op) 29 goto retry 30 return result </pre>	<pre> Input: int key, node_t* root Output: true: if the key is added to the set, false otherwise 1 node_t* parent 2 node_t* curr 3 node_t* new_node = NULL 4 operation_t* parent_op 5 operation_t* node_op 6 operation_t* cas_op 7 int result 8 while TRUE do 9 result = seek(key, &parent, &parent_op, &node, &node_op, root, root) 10 if result == FOUND && ! node→deleted then 11 return FALSE 12 if new_node == NULL then 13 new_node = alloc_node(k,v,0) 14 bool_t is_left = (result == NOT_FOUND_L) 15 node_t* old 16 if is_left then 17 old = node→left 18 else old = node→right 19 20 cas_op = alloc_op() 21 if result == FOUND && node→deleted then 22 cas_op→insert_op.is_update = TRUE 23 cas_op→insert_op.is_left = is_left 24 cas_op→insert_op.expected = old 25 cas_op→insert_op.new = new_node 26 if CAS(&node→op, node_op, FLAG(cas_op, INSERT)) == node_op then 27 help_insert(cas_op, node) 28 return TRUE </pre>
--	--

Algorithm 3: Seek

Algorithm 4: Insert

168 The *operation-descriptor* is either an, *insert_op* containing information for
 169 an intended insert, or a *rotate_op*: containing information for a rotation. An
 170 *insert_op* contains, *new*: pointer to the new node to be inserted, *expected*: pre-
 171 vious node, *isLeft*: indicating whether the new node would be inserted as a
 172 left-child or right-child, and *is_update*: indicating whether it is a new node to be
 173 inserted or an existing deleted node to be made part of the set again. A *rotate_op*
 174 descriptor contains in addition to the nodes involved in the corresponding opera-
 175 tion, a *state* field which can have three values, UNDECIDED: all nodes involved
 176 in the operation have not been grabbed yet, GRABBED: all the nodes that are

177 needed have been grabbed, ROTATED: indicating that the rotation has been
 178 completed.

179 As in [9, 17], we use two least significant bits of a pointer address to store
 180 the operation status. The operation pointer, *op*, of a node could be in one of
 181 following statuses, NONE: meaning no operation is currently intended, MARK:
 182 this node can be physically removed, ROTATE: a rotate operation is intended,
 183 and INSERT: an insert operation is intended. The following macros are used to
 184 modify *op*'s status, *FLAG(op, status)*: sets the operation pointer to one of the
 185 above status, *GETFLAG(op)*: returns the status of *op*, *UNFLAG(op)*: resets the
 186 status to NONE.

187 **Search:** The *search* algorithm, outlined in algorithm 2, is mostly similar to
 188 the sequential search with additional checks due to other concurrent inserts. It
 189 traverses from the *root*² to a leaf node looking for the *key*. If found, it breaks out
 190 of the loop as shown at lines 17-18. A further check is needed to see if the *key*
 191 is deleted from the set. In this case, if the corresponding node's *deleted* field is
 192 set and its *operation* field is flagged as INSERT then the new *key* which is to be
 193 inserted, is compared with the *key* that the *search* is looking for at lines 19-21. If
 194 a match is found then the algorithm returns *true* line 22. If *deleted* bit is not set
 195 and the node is found then the algorithm simply returns the true at line 24. If
 196 the *key* is not found the *search* algorithm returns false line 24. Synchronization
 197 of any form is not needed in this *search* algorithm. Also, the *search* algorithm
 198 never restarts.

199 **seek:** Algorithm 3 outlines the *seek* method which is used by both insert and
 200 delete algorithms to locate the position or potential position of a *key*, start-
 201 ing from a specified point *aux_root* in the tree. The position is returned in the
 202 variables pointed to by arguments *parent* and *node* and the values of their *op*-
 203 *eration* fields are returned in the variables pointed to by arguments *parent_op*
 204 and *node_op*. The result of the seek can be one of three values, FOUND: if the
 205 key was found, NOT_FOUND_L: if the *key* was not found but would be posi-
 206 tioned at the left child of the *node* if it were inserted, NOT_FOUND_R: similar
 207 to NOT_FOUND_L but for the right child. The variable *next* is used to point to
 208 the next node along the seek path, the *parent* and *parent_op* are used to record
 209 previous node. The check at lines 7-10 handles the case when the root has an
 210 ongoing operation to add to the empty tree. The seek loop traverses nodes one
 211 at a time until either the *key* is found or a null link is reached. Line 27 checks
 212 whether the node that is found does has an ongoing operation. If the node has
 213 an ongoing operation, the appropriate helping is done and seek restarts at lines
 214 28-29.

215 **Insert:** The *Insert* algorithm, algorithm 4, begins by calling the *seek* method
 216 at line 9. Depending on the result of the *seek* method : case 1: the *key* is found

² For implementation purpose, the first *node* to be inserted in the empty tree is kept
 as the right child of a fixed node *root* which has *key* assumed to be infinity

217 and the corresponding *node* is not *deleted*; case 2: the *key* is found and the
 218 corresponding *node* is *deleted*; case 3: the *key* is not found in the tree. For case
 219 1, the *insert* method returns false at line 11. For both cases 2 and 3, a *new-*
 220 *node* and an *insert_op operation-descriptor* are allocated with all the necessary
 221 information at lines 12-13 and 20-25 respectively. Furthermore, for case 2 the
 222 *is_update* field of *insert_op* descriptor is set to true at lines 21-22. Then a CAS
 223 is used to flag the *operation* field of the node to INSERT. If successful, the
 224 *help_Insert* method is called to complete the physical insertion, otherwise it
 225 retries. If the *help_insert* method finds *is_update* field set it simply sets off the
 226 *deleted* bit of the node to make it part of the set again otherwise it inserts the
 227 *new-node*.

228 **Delete:** Similar to the *Insert*, the *Delete* method starts by calling the *seek*
 229 method at line 6. If a *node* with the desired *key* is not found, it simply returns
 230 false at lines 7-8. If the *key* is found and is *deleted*, it further checks whether the
 231 node is flagged as INSERT, since an ongoing *insert* operation might be trying
 232 to insert the same *key*. If the *node* is not flagged as INSERT, the *delete* method
 233 returns false as the *node* is already deleted at lines 9-11. If the *key* is found and
 234 not deleted, the *delete* method tries to set the *node*'s *deleted* field to *true* and
 235 returns *true* on successful CAS at line 12-15.

236 **tree-maintenance:** The tree maintenance actions involve checking for balance
 237 condition violations, rotation operations, and physical removal of deleted nodes
 238 having at most one child. The latter two operations are started by a maintenance-
 239 thread and can be completed by any other thread³.

240 The maintenance-thread repeatedly performs a depth-first traversal in the
 241 background in which, at each node, it looks to see if the node can be removed
 242 and then adjusts heights. If it finds any node with *deleted* field set and at most
 243 one child, it then tries to flag the *operation* field of the node to MARK. If
 244 marking of the node is successful, the *help-marked* method outlined in algorithm
 245 7 is called to physically remove the node from the tree. After adjusting heights,
 246 a check is performed to see if a balance violation has occurred. At any node,
 247 the balance condition is said to be violated if *right-height* and *left-height* of the
 248 node differ by 2 or more. If there is a balance violation at any node further
 249 checks are performed to determine whether a single rotation (left or right) or
 250 double rotations are required. Once the type (left or right) and number (single
 251 or double) of rotation operations are determined the rotation process is started
 252 as listed in the *left_rotate*, algorithm 6, for the left rotation (similarly for the
 253 right rotation).

254 In the *left_rotate* method checks are performed to ensure that nodes that
 255 are involved in rotation operation are still intact at lines 2-6. The check at line
 256 7 is performed to see whether double rotations are needed. After this check, a
 257 *rotate_op* descriptor is allocated at line 10 using appropriate values.

³ Detailed explanations and full algorithm for tree-maintenance can be found in [13]


```

Input: int k, node_t* root
Output: true: if the node is found (not
          already deleted) and deleted
          from the set; else false

1 node_t* parent
2 node_t* node
3 operation_t* parent_op
4 operation_t* node_op
5 while TRUE do
6     int res = seek(k, &parent,
7                 &parent_op, &node, &node_op,
8                 root, root)
9     if (res ≠ FOUND) then
10        return FALSE
11    if node→deleted then
12        if GETFLAG(node→op) ≠
13            INSERT then
14                return FALSE
15    else
16        if GETFLAG(node→op) ==
17            NONE then
18            if CAS(&node→deleted,
19                FALSE, TRUE) ==
20                FALSE then
21                return TRUE

```

Algorithm 5: delete

```

259
int leftrotate(node_t* parent, int
dir, bool_t rotate)
1 node_t* node
2 if parent→removed then
3     return 0
4 node = (dir == 1) ?
5     (node_t*)parent→right :
6     (node_t*)parent→left
7 if ISNULL(node) ||
8     ISNULL(node→right) then
9     return 0
10 if (node→right→lh -node→right→rh)
11     > 0 && !rotate then
12     return 3; //double rotation
13 if GETFLAG(parent→op) == NONE
14     then
15     operation_t* rotateOp
16     rotateOp =
17     alloc_rotateop(parent,node,
18     node→right, parent→op,
19     node→op,
20     node→right→op,dir,FALSE)
21 if CAS(&(parent→op),
22     rotateOp→rotate_op.pOp,
23     FLAG(rotateOp,ROTATE)) ==
24     rotateOp→rotate_op.pOp then
25     help_rotate(rotateOp, parent,
26     node, node→right);
27 return 1;

```

Algorithm 6: left_rotate

```

void help(node_t* parent,
operation_t* parent_op, node_t*
node, operation_t* node_op)
1 if GETFLAG(node_op) == INSERT
2     then
3     help_insert(UNFLAG(node_op),
4     node)
5 else if GETFLAG(parent_op) ==
6     ROTATE then
7     help_rotate((UNFLAG(parent_op),
8     parent,node,
9     parent_op→rotate_op.child)
10    )
11 else if GETFLAG(node_op) ==
12     MARK then
13     help_marked(parent, parent_op,
14     node)

void help_insert(operation_t* op,
node_t* dest)
1 if op→insert_op.update then
2     CAS(&dest→deleted, TRUE,
3     FALSE)
4 else
5     node_t** address = NULL
6     if op→insert_op.is_left then
7         address = (node_t**)
8         &(dest→left)
9     else
10        address = (node_t**)
11        &(dest→right)
12    CAS(address,
13        op→insert_op.expected,
14        op→insert_op.new)
15    CAS(&(dest→op), FLAG(op,
16        INSERT), FLAG(op, NONE))

void help_marked(node_t* parent,
operation_t* parent_op, node_t*
node)
1 node_t* child, address
2 if ISNULL((node_t*) node→left) then
3     if ISNULL((node_t*)
4     node→right) then
5         child =
6         (node_t*)SETNULL(node)
7     else
8         child = (node_t*) node→right
9 else
10    child = (node_t*) node→left
11 node→removed = TRUE
12 operation_t* cas_op = alloc_op()
13 cas_op→insert_op.is_left = (node ==
14     parent→left)
15 cas_op→insert_op.expected = node
16 cas_op→insert_op.new = child
17 if CAS(&(parent→op), parent_op,
18     FLAG(cas_op, INSERT)) ==
19     parent_op then
20     help_insert(cas_op, parent)

```

Algorithm 7: help,help_insert,help_marked

261 Then, an attempt is made to insert *rotate_op* to the *parent* of the *node*
 262 where the violation has occurred using CAS. If successful, a call is made to
 263 the *help_rotate* method.

```

Input: operation_t* op, node_t* parent, node_t* node, node_t* child
1 int seen_state = op->rotate_op.state
2 retry:
3 if seen_state == UNDECIDED then
4   if GETFLAG(node->op) == NONE ||
   GETFLAG(node->op) == ROTATE
   then
5     if GETFLAG(node->op) == NONE
   then
6       CAS(&(op->rotate_op.node->op),
   node->op, FLAG(op, ROTATE))
7     if GETFLAG(node->op) ==
   ROTATE then
8       nodegrabbed:
9       if GETFLAG(child->op)
   == NONE ||
   GETFLAG(child->op)
   == ROTATE then
10      if GETFLAG(child->op)
   == NONE then
11      CAS(&(op->rotate_op.
   child->op), child->op,
   FLAG(op, ROTATE))
12      if GETFLAG(child->op)
   == ROTATE then
13      CAS(&
   (op->rotate_op.state),
   UNDECIDED,
   GRABBED)
14      seen_state = GRABBED
15      else goto nodegrabbed
16
17      else
18      help(node, node->op, child,
   child->op)
19      goto nodegrabbed
20
21      else
22      goto retry
23
24      else
25      help(parent, parent->op, node,
   node->op)
   goto retry
26 if seen_state == GRABBED then
27   if op->rotate_op.rightR then
28     // right rotation
29     //allocate newnode with
30     appropriate children and heights
31     //carry out rotation steps
32     using CAS
33   else
34     // left rotation
35     //allocate newnode with
36     appropriate children and heights
37     //carry out rotation steps
38     using CAS
39   //parent pointer swing
40   if op->rotate_op.rightR then
41     CAS(&op->rotate_op.parent->left,
   op->rotate_op.node, child)
42   else
43     CAS(&op->rotate_op.parent
   ->right, op->rotate_op.node, child)
44   //adjust child and parent heights
45   CAS(&(op->rotate_op.state),
   GRABBED, ROTATED)
46   seen_state = ROTATED;
47
48   if seen_state == ROTATED then
49     //Clear parent, node, child operation
50     fields from ROTATE to NONE
51
52 bool_t result = (seen_state == ROTATED)
53 return result;

```

Algorithm 8: help_rotate

264 The *help_rotate* method, outlined in algorithm 8 can be called by any thread
 265 if it finds that the *operation* field of a node is flagged as ROTATE. Any thread
 266 executing *help_rotate* tries to the grab remaining nodes namely, the *node* and
 267 the *child* by flagging their *operation* fields to ROTATE at lines 2-24. Failing to
 268 flag any node means that there is an ongoing operation. In this case, the thread
 269 helps the ongoing operation first, lines 18 and 23. Once both *node* and *child* are
 270 flagged the *state* field of *rotate_op* descriptor is updated to a value GRABBED at
 271 line 13. The rest of the operation is carried out (lines 25-40) as shown in figure 2.

272 It should be noted that all of the methods involved in tree balancing except the
 273 *help_rotate* are not exposed to any thread other than the maintenance-thread.

274 **help_insert/Marked:** Both these methods are called only after inserting ap-
 275 propriate flags to the *operation* field of the *node*. Once a *node* is flagged it will
 276 never be un-flagged until the operation is completed. The former adds a new
 277 node as a child of the node that was flagged as INSERT or simply sets off the
 278 *deleted* field depending on the *is_update* field of the node, and the later physically
 279 removes the marked node from the tree.

280 3.3 Correctness

281 **Lock-freedom :** The non-blocking property of the algorithm is explained by
 282 describing the interactions that can occur between the threads that are executing
 283 read and write operations. A *search* will either locate the *key* it is searching for
 284 or terminate at a leaf node. The search operation never restarts. An *Insert*
 285 or *Delete* operation retries until it gets clean nodes through the *seek* method.
 286 Then, the *insert* tries to flag node as INSERT which cannot be undone until the
 287 physical insertion is completed. The *Delete* operation will only set the *deleted*
 288 bit when there is no other *insert* operation going on concurrently. The *seek*
 289 method restarts when it finds the leaf node has an ongoing operation when
 290 called from *Insert*, or if it finds the node having the *key* to be deleted has an
 291 ongoing operation when called from *Delete*. In both cases if the *operation* field
 292 of the node contains INSERT, ROTATE or MARK, this means an operation
 293 has been applied to the tree so system-wide progress was achieved. Assuming
 294 that there is no infinite addition of nodes on its search path, the *seek* method
 295 will eventually return clean nodes. An *Insert* operation could also restart if
 296 it fails to flag the node returned by *seek* to INSERT. In this case also it has
 297 encountered an ongoing operation. Any thread executing an operation will help
 298 the ongoing operation to its completion before restarting its own operation.
 299 Similarly, a rotation starts only when the maintenance-thread successfully flags
 300 the parent of the node where balance violation has occurred. If flagging of other
 301 nodes involved rotation fails in the *help_rotate* method, this means there is an
 302 ongoing operation and the process completes that first before coming back to
 303 flagging the node to ROTATE.

304 **Linearisability :** To prove the linearisability of RBLFBST, we first define
 305 the linearisation points of the *Search*, *Insert*, and *Delete* operations. The *search*
 306 operation can have two possible outcomes: a *key* is found in the set or not. The
 307 linearisation point for finding a *key* is the point at which an atomic read of the
 308 *key* has occurred at line 11. As our design allows a deleted *key* to be present in
 309 the tree it has to pass check at line 19. If the search does not find the *key*, it
 310 will linearise reading a NULL link at either line 13 or 15. If the tree is empty
 311 the search will linearise at reading the null link at line 7. A successful *Insert*
 312 operation will linearise when the *operation-descriptor* is successfully inserted to

313 the *node* returned by the seek algorithm at line 26. Failure of insert operation
 314 will have linearisation point at line 17 of the *seek* algorithm where it reads the
 315 *key* to be inserted already present in the tree. However, it has to verify if the
 316 *key* is logically deleted by failing the test at line 11 of the *Insert* algorithm.
 317 Similarly, the successful *Delete* operation will linearise at the successful CAS at
 318 line 14 of the delete algorithm. Linearisation point of the failed *Delete* will occur
 319 at line 7 of the delete algorithm where it is verified that the *key* is not present in
 320 the tree or is logically deleted by passing the checks at lines 10-12. An elaborate
 321 correctness discussion can be found in [13].

322 4 Experimental Results

323 To evaluate performance of RBLFBST we compared it with following recently
 324 published implementations : (i) non-blocking internal BST denoted by bst-howley
 325 [17]. (ii) lock-free external binary BST denoted by , bst-Aravind [14]. (iii) lock-
 326 based partially internal BST with balancing operations denoted by bst-bronson
 327 [15]. (iv) lock-based internal BST with balancing operations denoted by bst-
 drachsler [5].

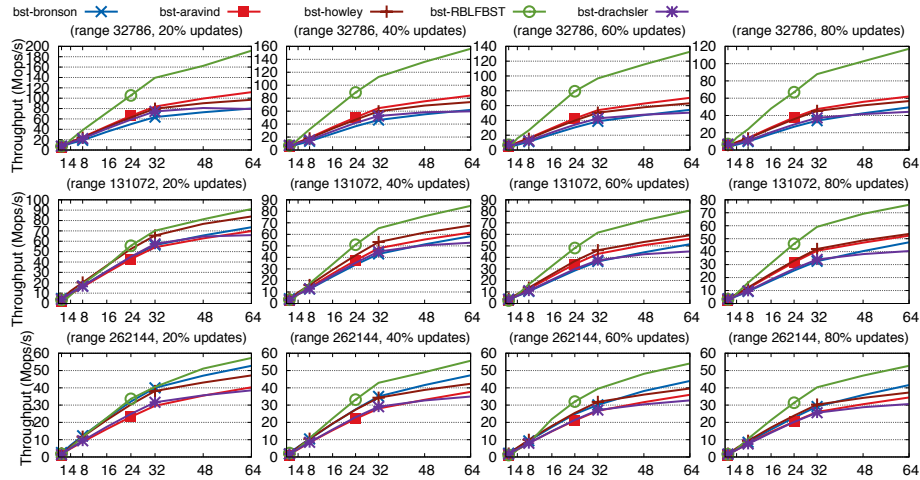


Fig. 3: Throughput in Mops (million of operation per second) for different algorithms for varying parameters. Each row shows results for variations in the distribution of operations. Each column shows results for different key ranges.

328 All experiments were conducted on a machine with 32 core AMD Ryzen
 329 Threadripper 2990WX processor, 64 hardware threads, and 32 GB RAM running
 330 x86_64 version of Ubuntu 18.04. All codes were implemented in C and compiled
 331 using gcc version 7.5.0 using optimization level O3. The source codes for other
 332 implementations were obtained from ASCYLIB [8].

334 The comparison of RBLFBST with other implementation was done varying
 335 three parameters: the key range ($[0, 32786]$, $[0, 131072]$ and $[0, 262144]$), the
 336 distribution of operations, and the number of threads concurrently executing

337 on the data structures (1 to 64). Operations distribution considered were: (a)
 338 Write dominated workload: 80% updates, 20% search (b) mixed workload: 60%
 339 updates, 40% search (c) mixed workload: 40% updates, 60% search (d) read
 340 dominated workload: 20% updates, 80% search. Update operations had an equal
 341 number of insert and delete operations. All the operation types and the keys
 342 were generated using a pseudo-random number generator which were seeded
 343 with different values for each run. Each run was carried out for 5 seconds, and
 344 the results were collected averaging throughput over 20 runs. To mimic steady
 345 state the tree was filled with half the key range for each run.

346 The graphs in figure 3 show the comparisons of other implementations against
 347 RBLFBST. Overall, RBLFBST scales very well as the number of threads in-
 348 creases. For smaller key range 2^{15} , RBLFBST outperforms its nearest competitor
 349 bst-aravind by a maximum of 90% and 85% in read-dominated (80% and 60%
 350 search respectively), by 87% and 70% (60% and 80% updates respectively) in
 351 write-heavy workloads. Our tree beats the other two lock-based trees with bal-
 352 ancing operations (bst-drachsler and bst-bronson) by 142% in 80% search work-
 353 load and more than 85% in update heavy workload for the same key range. The
 354 performance of RBLFBST for 2^{17} key range again is better by 43% and 37% than
 355 its nearest lock-free competitor bst-howley in 60% and 80% update workload re-
 356 spectively. Similarly, it performs 56% and 61% better than the closest lock-based
 357 competitor bst-bronson in 60% and 80% update workload respectively. For the
 358 key range 2^{18} , RBLFBST beats bst-bronson by 24-38% and 20-30%, bst-howley
 359 by 10-31% and 20-30%, bst-aravind by 11-53% and 47-50% and bst-drachsler
 360 by 11-72% and 2-65% for the workload containing 80% and 60% updates re-
 361 spectively. The better performance of RBLFBST is due to the fact that it uses
 362 less number of expensive operations (CAS) than other lock-free implementations
 363 for insert and delete operations combined, thereby allowing more concurrency.
 364 The performance of our design goes down relatively for read-heavy workloads
 365 (80% read) particularly for the larger key ranges (2^{18} , 2^{17}). Larger key range will
 366 grow tree longer and more rotations will be performed. This is when the effect
 367 of threads helping rotation operation which can use up to 10 CAS instructions
 368 is clearly visible. However, the performance of RBLFBST in such cases is still
 369 better than all other implementations.

370 5 Conclusion and Future work

371 In this work, we presented a relaxed balanced lock-free binary search tree. In our
 372 design, all the set operations are lock-free. The search operation in our algorithm
 373 is oblivious to any structural changes done by other operations and is also free
 374 of any additional synchronization. Our results show its concurrent performance
 375 to be very good compared with other concurrent BSTs. We have discussed the
 376 correctness of RBLFBST operations. We are working on formal verification of
 377 our algorithm. We are also planning to apply the separation of tree balancing
 378 operations as well as relaxing delete operation to other lock-free balanced binary
 379 trees designs as future work.

380 **References**

- 381 1. Afek, Y., Kaplan, H., Korenfeld, B., Morrison, A., Tarjan, R.E.: Cbtree: A practical
382 concurrent self-adjusting search tree. In: Aguilera, M.K. (ed.) Distributed
383 Computing. pp. 1–15. Springer Berlin Heidelberg, Berlin, Heidelberg (2012)
- 384 2. Barnes, G.: A method for implementing lock-free shared-data structures. In: Proceedings
385 of the Fifth Annual ACM Symposium on Parallel Algorithms and Architectures. pp. 261–270. SPAA '93 (1993)
- 386 3. Bougé, L., Vallés, J., PeyPOCH, X.M., Schabanel, N.: Height-relaxed avl rebalancing:
387 a unified, fine-grained approach to concurrent dictionaries (1998)
- 388 4. Crain, T., Gramoli, V., Raynal, M.: A contention-friendly binary search tree. In:
389 Wolf, F., Mohr, B., and Mey, D. (eds.) Euro-Par 2013 Parallel Processing. pp. 229–
390 240. Springer Berlin Heidelberg, Berlin, Heidelberg (2013)
- 391 5. Drachsler, D., Vechev, M., Yahav, E.: Practical concurrent binary search trees via
392 logical ordering. In: Proceedings of the 19th ACM SIGPLAN Symposium on Principles
393 and Practice of Parallel Programming. pp. 343–356. PPOPP '14, Association
394 for Computing Machinery, New York, NY, USA (2014)
- 395 6. Faith Ellen, Panagiota Fatourou, E.R., van Breugel, F.: Non-blocking binary search
396 trees. In Proceedings of the 29th ACM SIGACT-SIGOPS symposium on Principles
397 of distributed computing, PODC '10 pp. 131–140 (2010)
- 398 7. Fraser, K.: Practical Lock freedom. Ph.D. thesis, King's College ,University of
399 Cambridge (September 2003)
- 400 8. Guerraoui, R.: Ascylib (July 2020), <https://dcl.epfl.ch/site/optik>
- 401 9. Harris, T.L.: A pragmatic implementation of non-blocking linked-lists. In: Proceedings
402 of the 15th International Conference on Distributed Computing. pp. 300–314.
403 DISC '01, Springer-Verlag, Berlin, Heidelberg (2001)
- 404 10. Herlihy, M., Shavit, N.: The Art of Multiprocessor Programming. Morgan Kauf-
405 mann Publishers Inc., San Francisco, CA, USA (2008)
- 406 11. Kessels, J.L.W.: On-the-fly optimization of data structures. *Commun. ACM*
407 **26**(11), 895–901 (Nov 1983)
- 408 12. Larsen, K.S.: Avl trees with relaxed balance. *Journal of Computer and System*
409 *Sciences* **61**(3), 508–522 (2000)
- 410 13. Manish Singh, Lindsay Groves, A.P.: A relaxed balanced lock-free binary search
411 tree. Tech. rep. (2020), <https://ecs.wgtn.ac.nz/Main/TechnicalReportSeries>
- 412 14. Natarajan, A., Mittal, N.: Fast concurrent lock-free binary search trees. In: Proceedings
413 of the 19th ACM SIGPLAN Symposium on Principles and Practice of
414 Parallel Programming. pp. 317–328. PPOPP '14, New York, NY, USA (2014)
- 415 15. Nathan G. Bronson, Jared Casper, H.C., Olukotun, K.: A practical concurrent
416 binary search tree. ACM SIGPLAN Symposium on Principals and Practice of
417 Parallel Programming (2010)
- 418 16. Nurmi, O., Soisalon-Soininen, E.: Uncoupling updating and rebalancing in chromatic
419 binary search trees. In: Proceedings of the Tenth ACM SIGACT-SIGMOD-
420 SIGART. pp. 192–198. PODS '91, New York, NY, USA (1991)
- 421 17. Shane V. Howley, J.J.: A non blocking internal binary tree. SPAA (June 2012)
- 422 18. T. Brown, F. Ellen, E.R.: A general technique for non-blocking trees. In Proceedings
423 of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel
424 Programming (PPOPP) (2014)
- 425 19. Valois, J.D.: Lock-Free Data Structures. Ph.D. thesis, Rensselaer Polytechnic In-
426 stitute, Troy, NY, USA (1996)
- 427