

On Concurrent Lock-Free Binary Search Trees

Arunmoezhi Ramachandran, Neeraj Mittal

Department of Computer Science, The University of Texas at Dallas, Richardson,
Texas 75080, USA

(arunmoezhi,neerajm)@utdallas.edu

Keywords: Concurrent Data Structure, Lock-Free Algorithm, Binary Search Tree

1 Introduction

A Binary Search Tree (BST) implements the dictionary abstract data type. In a BST, all the values in a node's left subtree are less than the node value and all the values in the node's right subtree are greater than the node value. Duplicates are not allowed. A BST supports three main operations, viz.: search, insert and delete. $\text{Search}(x)$ determines if x is present in the tree. $\text{Insert}(x)$ adds key x to the tree if it is not already present. $\text{Delete}(x)$ removes the key x from the tree if it is present.

Several algorithms have been proposed for non-blocking binary search trees. Ellen *et al.* proposed the first practical lock-free algorithm for a concurrent binary search tree in [1]. Their algorithm uses an external (or leaf-oriented) search tree in which only the leaf nodes store the actual keys; keys stored at internal nodes are used for routing purposes only. Howley and Jones have proposed another lock-free algorithm for a concurrent binary search tree in [2] which uses an internal search tree in which both the leaf nodes as well as the internal nodes store the actual keys. One advantage of using an internal search tree is that its memory footprint is half that of an external search tree. For large size trees, search time dominates and our experimental results shows that internal BST tend to perform better than external BSTs. Natarajan and Mittal have proposed another lock-free external BST in [3]. The key change in this algorithm is that it operates at edge level and the ones described by Ellen *et al.* [1] and Howley and Jones [2] operates on node level. A node (or edge) level operation mean that nodes (or edges) are marked for deletion. Operating at edge level blocks fewer operations than operating at node level. So edge level operation provides more concurrency when there is high contention. Hence the algorithm described by Natarajan and Mittal [3] performs well for smaller trees with high contention. We have designed our algorithm to get the benefits of both the worlds. Our algorithm like the one proposed by Howley and Jones [2] is an internal BST and like the one proposed by Natarajan and Mittal [3] operate on edge level.

2 Lock-Free Algorithm

Every operation in our algorithm begins with a seek phase.

Algorithm 1: structures used

```

1 struct Node{
2     {Boolean,key} markAndKey;
3     {Boolean,Boolean,Boolean,NodePtr} child[2];
4     Boolean readyToReplace;
5 };
6 struct seekRecord{
7     NodePtr node; NodePtr parent;
8     NodePtr lastUParent; NodePtr lastUNode;
9     NodePtr injectionPoint;
10 };
11 struct State{
12     NodePtr node; NodePtr parent;
13     Key key;
14     enum mode{ INJECTION, DISCOVERY, CLEANUP };
15     enum type{ SIMPLE, COMPLEX };
16     seekRecPtr seekRec;
17 };

```

2.1 Overview

Seek: The operation traverses the search tree from the root node until it finds the target key or it reaches a leaf node. We refer to the path traversed by the operation in the seek phase as the access path. The operation then compares the target key with the value stored in node. Depending on the result of the comparison and the type of the operation, the operation either terminates or moves to the next phase. For certain cases where some key may have moved up the tree, the seek operation may have to restart (will be discussed later). We now describe the next steps for each of the type of operation one-by-one.

Search: A search operation invokes seek and returns *true* if the stored key matches the target key; otherwise it returns *false*.

Insert: An insert operation invokes seek and returns *false* if the stored key matches the target key; otherwise it moves to the execution phase. In the execution phase, it attempts to insert the key into the tree as a child node of the last node in the access-path using a CAS instruction. If the instruction succeeds, then the operation returns *true*; otherwise, it restarts from the seek phase after possibly helping.

Delete: A delete operation invokes seek and returns *false* if the stored key does not match the target key; otherwise it moves to the execution phase. In the execution phase, it attempts to remove the key stored in the last node in the access path. let the last node be referred to as *target node*. There are two cases

Algorithm 2: $\text{seek}(key, \text{seekRec})$

```

18 while true do
    // create two local seek records: cSeek (current seek record) and
    // pSeek (previous seek record) used for the traversal
19 while true do
20      $\langle *, cKey \rangle := \text{curr} \rightarrow \text{markAndKey}$ ; // key in the curr of cSeek
21     if  $key = cKey$  then // key found; stop the traversal
22         done := true; break;
23     which :=  $key < cKey$  ? LEFT: RIGHT;
24      $\langle n, d, p, \text{address} \rangle := \text{curr} \rightarrow \text{child}[\text{which}]$ ; // read the next edge
25     if n then // null flag is set; reached a leaf node
26         if key stored in anchorNode has not changed then
27             done := true; break; // use data from cSeek
28         else if anchorNode of cSeek & pSeek matches then
29             done := true; break; // use data from pSeek
30         else
31             break; // after copying cSeek to pSeek
32     if which = RIGHT then // next edge to be traversed is a right edge
33         anchorNode := curr; // keep track of curr node
34         anchorKey := cKey; // and its key
35     prev := curr; curr := address; // traverse the next edge
36     if not (d or p) then // keep track the last unmarked edge
37         lastUParent := prev; lastUNode := curr;
38 if done then
    // initialize the appropriate seek record (cSeek or pSeek)
39 return;

```

depending on if the *target node* is a binary node (has two children) or not (at most one child). In the first case, the operation is referred to as complex delete. In the second case, it is referred to as simple delete. In the case of simple delete (as shown in Fig. 1), the target node is removed by changing the pointer at the parent of the *target node*. In the case of complex delete (as shown in Fig. 2), the key to be deleted is replaced with the next largest key in the tree (which will be stored in the leftmost node of the right subtree of the *target node*).

2.2 Details of the Algorithm

We use sentinel keys and nodes to handle the boundary cases easily. A tree node in our algorithm consists of three fields: (i) *markAndKey* which contains the key stored in the node, (ii) *childarray*, which contains the addresses of the left and right children and (iii) *readyToReplace*, which is a boolean flag used by complex delete operation to indicate if a node can be replaced with a fresh copy of it.

Algorithm 3: Search(*key*)

```

40 seek( key, mySeekRec);
41  $\langle *, nKey \rangle := mySeekRec \rightarrow node \rightarrow markAndKey$ ;
42 if key = nKey then return true;
43 else return false;

```

Algorithm 4: Insert(*key*)

```

44 while true do
45   seek( key, mySeekRec);
46    $\langle *, nKey \rangle := mySeekRec \rightarrow node \rightarrow markAndKey$ ;
47   if key = nKey then return false;
48   newNode := create a new node and initialize its fields;
49   which := key < nKey ? LEFT: RIGHT;
50    $\langle *, *, *, address \rangle := mySeekRec \rightarrow injectionPoint$ ;
51   out := CAS(node → child[which],  $\langle 1, 0, 0, address \rangle$ ,  $\langle 0, 0, 0, newNode \rangle$ );
52   if out then return true;
53    $\langle *, d, p, address \rangle := node \rightarrow child[which]$ ; // find out why the CAS failed
54   if not (d or p) then continue; // CAS failed due to another insert op
55   deepHelp(mySeekRec → lastUNode, mySeekRec → lastUParent);

```

This algorithm like the algorithm described by Natarajan and Mittal [3], operates on edge level. A delete operation obtains ownership of the edges it needs to work on by marking them. To enable marking we steal two bits from the child addresses of a node. To avoid ABA problem, as in Howley and Jones [2], we steal another bit from the child address referred to as *nullFlag* and use it to indicate whether the address points to a null or a non-null value. So when an address changes from a non-null value to a null value, we only set the *nullFlag* and the contents of the address are not modified. As complex delete replaces a key in a node being deleted, a flag is required to identify if the key in a node has changed. So we steal a bit from the key field and use it as a mark bit. If the mark bit is set, it denotes that the key in the node has changed.

We next describe the details of the seek phase, which is executed by all operations (search as well as modify) after which we describe the details of the execution phases of insert and delete operations.

The Seek Phase

A seek operation keeps track of the node in the access path at which the last "right turn" was taken (i.e., right edge was followed). Let us call this node as *anchorNode*. Upon reaching the last node, it compares the stored key with the target key. If they do not match, then it is possible that the key may have moved up in the tree. So key stored in the *anchorNode* is checked to see if has changed. If the key has changed then the seek operation restarts. If the key

Algorithm 5: Inject(*state*)

```

56 node := state → node // try to set the delete flag on the left edge
57 while true do
58    $\langle n, d, p, left \rangle := node \rightarrow child[LEFT]$ ;
59   if d or p then return false; // edge is already marked
60   out := CAS(node → child[LEFT],  $\langle n, 0, 0, left \rangle$ ,  $\langle n, 1, 0, left \rangle$ );
61   if out then break; // retry from beginning of while loop
62   updateModeAndType(state) // mark right edge, update mode and type
63   return true;

```

Algorithm 6: updateModeAndType(*state*)

```

64 node := state → node // retrieve the address from the state record
65 if node → child[RIGHT] ≠  $\langle *, 1, *, * \rangle$  then // mark right edge if unmarked
66   BTS(node → child[RIGHT], DELETE_FLAG);
67  $\langle m, * \rangle := node \rightarrow markAndKey$ ;
68  $\langle lN, *, *, * \rangle := node \rightarrow child[LEFT]$ ;  $\langle rN, *, *, * \rangle := node \rightarrow child[RIGHT]$ ;
69 if lN or rN then // update the op mode and type
70   if m then
71     state → type := COMPLEX; node → readyToReplace := true;
72   else
73     state → type := SIMPLE; state → mode := CLEANUP;
74 else
75   state → type := COMPLEX;
76   if readyToReplace then state → mode := CLEANUP;
77   else state → mode := DISCOVERY;
78 return ;

```

has not changed, then the key stored in the *anchorNode* is checked to see if it is undergoing deletion (by checking if its left child edge is marked). If the key is not undergoing deletion, then the seek operation terminates. If the key is undergoing deletion, then it checks if the *anchorNode* of the current traversal matches with the *anchorNode* of the previous traversal. If they match, then the seek operation terminates by returning the results of the previous traversal; otherwise it restarts.

In the case of insert operation, the seek operation also returns the injection point which is the current contents of the child location at which the insert will occur.

Execution Phase of an Insert Operation

In the execution phase, an insert operation creates a new node containing the target key. It then adds the new node to the tree at the injection point using a CAS instruction. If the CAS succeeds, then (the new node becomes a part

Algorithm 7: findSmallest(*node, right, seekRec*)

```

// find the smallest key in the subtree rooted at the right child
79 lastUParent := node; lastUNode := right; prev := node; curr := right;
80 while true do
81    $\langle n, d, p, left \rangle := curr \rightarrow child[LEFT]$ ;
82   if  $n$  then break;
83   prev := curr; curr := left; // traverse the next edge
84   if not ( $d$  or  $p$ ) then // keep track of the last unmarked edge
85     lastUParent := prev; lastUNode := curr
      // update the seek record
86   return;

```

of the tree and) the operation terminates; otherwise, the operation determines if it failed because of a conflicting delete operation in progress. If there is no conflicting delete operation in progress then the operation restarts from the seek phase; otherwise it performs helping (which will be described later) and then restarts from the seek phase.

Execution Phase of a Delete Operation

The execution of a delete operation starts in the *injection* mode. Once the operation has been injected into the tree, then it advances to the *cleanup* mode.

Injection Mode: In the injection mode, the delete operation marks the left child edge of the target node using a CAS instruction. If the CAS instruction succeeds, then the delete operation has been injected into the tree and is guaranteed to complete. Then the operation moves on to the *cleanup* mode. But if the CAS instruction fails, the operation performs helping and restarts from the seek phase (and stays in the injection mode).

Cleanup Mode: In the *cleanup* mode, the operation begins by marking the right child edge of the target node using a BTS (Bit Test and Set) instruction (this can also be done using a CAS instruction as well). Note that we maintain an invariant that edges which are once marked cannot be unmarked. Eventually the node is either removed from the tree (by simple delete) or replaced with a "new" node containing the next largest key (by complex delete). Further, a marked edge is changed only under a specific situation by a delete operation as part of helping (described later).

After marking both the edges, an operation checks whether the node is a binary node or not. If it is a binary node, then the delete operation is classified as a complex delete; otherwise it is classified as a simple delete. Let **A** be the target node, **B** be its parent node and **C** and **D** be the left and right child node respectively, of **A**. We now discuss the two types of delete operation.

Algorithm 8: Delete(*key*)

```

// initialize the state record
87 myState→mode:= INJECTION; myState→key:= key;
88 while true do
89   seek( key, mySeekRec);
90   node:= mySeekRec→node; parent:= mySeekRec→parent;
91   ⟨*, nKey⟩ := node→markAndKey;
92   if myState→key≠ nKey then
93     // the key does not exist in the tree
94     if myState→mode= INJECTION then return false;
95     else return true;
96   needToHelp:=false;
97   // perform appropriate action depending on the mode
98   if myState→mode= INJECTION then
99     myState→node:= node // store a reference to the node
100    out:= inject(myState) // attempt to inject
101    if not out then needToHelp:= true;
102    // mode would have changed if the op was injected
103    if myState→mode≠ INJECTION then
104      // if the node found by seek is different from the one stored
105      // in state record, then the node is already deleted
106      if myState→node≠ node then return true;
107      myState→parent:= parent // update parent with recent seek
108    if myState→mode= DISCOVERY then
109      findAndMarkSuccessor(myState);
110    if myState→mode= DISCOVERY then
111      removeSuccessor(myState);
112    if myState→mode= CLEANUP then
113      out:= cleanup(myState,0);
114      if out then return true;
115      else
116        ⟨*, nKey⟩ := node→markAndKey; myState→key:= nKey;
117        // help if helpee node is not the node of interest
118        if mySeekRec→lastUNode≠ node then needToHelp:=true;
119    if needToHelp then
120      deepHelp(mySeekRec→lastUNode,mySeekRec→lastUParent) ;

```

Algorithm 9: cleanup(*state*, *dFlg*)

```

    // retrieve the addresses from the state record
114 pWhich := edge of the parent which needs to be switched;
115 if state → type = COMPLEX then
116     newNode := a new copy of the node in which all the fields are unmarked;
    // try to switch the edge at the parent
117     out := CAS(parent → child[pWhich], ⟨0, dFlg, 0, node⟩, ⟨0, dFlg, 0, newNode⟩);
118 else
119     nWhich := non-Null child of the node being deleted;
120     ⟨n, *, *, address⟩ := node → child[nWhich];
121     if n then // set only the null flag; do not change the address
122         out := CAS(parent → child[pWhich], ⟨0, dFlg, 0, node⟩, ⟨1, dFlg, 0, node⟩);
123     else // change the address here by switching the pointer
124         out := CAS(parent → child[pWhich], ⟨0, dFlg, 0, node⟩, ⟨0, dFlg, 0, address⟩);
125 return out;

```

Algorithm 10: findAndMarkSuccessor(*key*, *seekRec*)

```

    // retrieve the addresses from the state record
126 node := state → node; seekRec := state → seekRec;
127 while true do
128     right := address of the right child;
129     findSmallest(node, right, seekRec);
130     succNode := seekRec → node; // retrieve succ node from seek record
131     left := address of the left child of the succNode;
    // try to set the promote flag & copy the node address on the
    // left edge using CAS
132     out := CAS(succNode → child[LEFT], ⟨1, 0, 0, left⟩, ⟨1, 0, 1, node⟩);
133     if out then break; // promote flag set; promotion will eventually succeed
    // reread the edge to see why the attempt to mark the edge failed
134     ⟨n, d, p, left⟩ := succNode → child[LEFT];
135     if p then
136         if left = node then
137             break // successor node has already been selected
138         else // the node found is a successor node for another delete operation
139             node → readyToReplace := true
140             if not n then continue; // the node found has since gained a left child
141             if d then the node found is undergoing deletion. So invoke helping;
142 updateModeAndType(state); // update the operation mode and type
143 return;

```

Algorithm 11: removeSuccessor(*state*)

```

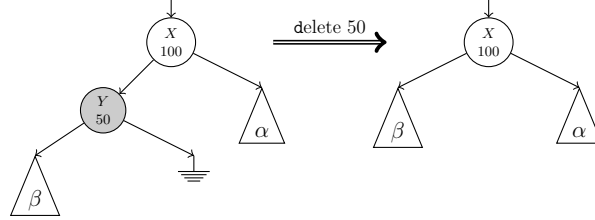
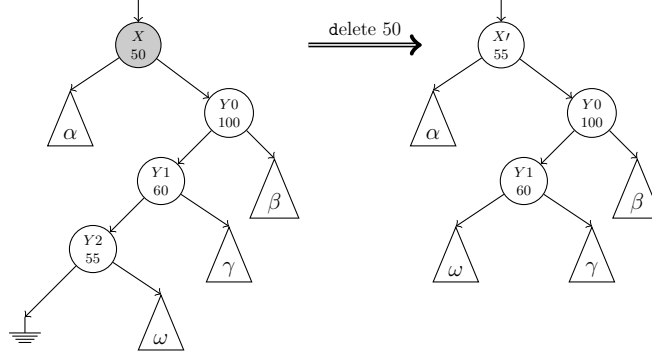
    // retrieve the addresses from the state record
144 node := state → node; seekRec := state → seekRec;
145 succNode := seekRec → node;
146 if promote flag not set on right child edge of succNode then
147   | BTS(succNode → child[RIGHT], PROMOTE_FLAG);
148 node → markAndKey := ⟨1, succNode → markAndKey⟩; // promote the key
149 while true do
150   | succParent := seekRec → parent; // retrieve parent of the succNode
151   | right := right child address of succNode;
152   | out := CAS(succParent → child[LEFT], ⟨0, 0, 0, succNode⟩, ⟨0, 0, 0, right⟩);
153   | if out then break; // successor removed successfully
    // invoke helping if needed
154   | findSmallest(node, right, seekRec);
155   | if seekRec → node ≠ succNode then break; // successor already removed
156 node → readyToReplace := true;
157 if state → parent ≠ null then updateModeAndType(state);
158 return;

```

Simple Delete: In this case, either C or D is a null node. Note that both C and D may be null nodes in which case A will be a leaf node. Without loss of generality, assume that D is a null node. Delete operation attempts to change the pointer at B that is pointing to A to point to C using a CAS instruction. If the CAS instruction succeeds, then the operation terminates; otherwise, the delete operation performs another seek operation. If the seek operation either fails to find the target key or return a target node different from A , then A has been already removed from the tree (by another operation as part of helping) and the operation terminates; otherwise, it attempts to remove A from the tree again. Note that the new seek operation may return a different parent node. This process may be repeated multiple times.

Complex Delete: In this case, both C and D are non-null nodes. The operation now performs the following steps:

1. Locate the next largest key in the tree, which is the smallest key in the subtree rooted at the right child of the target node A . We refer to this key as the successor key and the node storing this key as the successor node. Let S denote the successor node and T denote its parent node.
2. Claim the successor node. This involves marking both the child edges of S . Note that the left edge of S will be null. To distinguish between marking a target node (for deletion) and marking a successor node (for promotion), we steal two bits from the address and refer to them as *deleteFlag* and *promoteFlag*. The left edge of S is marked (i.e., *promoteFlag* is set) using a CAS instruction. As part of marking the left edge, we also store the address

Fig. 1. An illustration of a simple delete operation.**Fig. 2.** An illustration of a complex delete operation.

of the target node A in the left edge. This is done to enable helping in case the successor node is obstructing the progress of another operation. In case if the CAS instruction fails, the operation repeats from step 1. The right edge of S is marked using a BTS instruction.

3. Promote the successor Key. The successor key is copied into the target node. At the same time, the mark bit in the key is set to indicate that the key currently stored in the target node is the successor key and not the target key.
4. The successor node S is deleted by changing the pointer at T that is pointing to S to point to the right child of S using a CAS instruction. If the CAS instruction fails, then the operation performs helping if needed. It then finds the successor node again by performing a traversal starting from the right child of the target node A and repeats step 4. If the successor node is not found in the traversal, then it has been already removed from the tree (by another operation as part of helping) and the operation moves to step 5.
5. Note that, at this point, the original key in the target node has been replaced with the successor key. Further, the key as well as both its edges are marked. The target node is now replaced with a new node whose contents are same as that of the target node except that all the fields are unmarked. The target node is then replaced with a new node using a CAS instruction at the parent node. If the CAS instruction succeeds, then the operation terminates; otherwise, as in the case of simple delete, it performs another seek operation,

this time looking for the successor key. If the seek operation either fails to find the successor key or returns a target node different from A , then A has been already replaced (by another operation as part of helping) and the operation terminates. On the other hand if the seek operation finds the target node A with the successor key, it attempts to replace A again. Note that the new seek operation may return a different parent node. This process may be repeated multiple times.

Helping

To enable helping, whenever traversing the tree to locate either a target key or a successor key, we keep track of the last unmarked edge encountered in the traversal. Whenever an operation fails while executing a CAS instruction, it helps the operation in progress at the child end of the unmarked edge if different from its own operation. Note that, when traversing the tree looking for a target key, the last unmarked edge will always be found because of the sentinel keys. However, when traversing the tree looking for a successor key, the last unmarked edge may not always exist since the traversal starts from the middle of the tree from the right child of a target node. Recall that A denotes a target node and D its right child node. If no unmarked edge is found during the traversal from D , then helping is performed along the edge (A, D) under the following conditions (i) delete operation at D is of type simple delete, or (ii) delete operation at D is of type complex delete and is currently at its last step (replacing the target node with a new node).

3 Experimental Evaluation

Experimental Setup: We conducted our experiments on an X86_64 AMD Opteron 6276 machine running GNU/Linux operating system. We used gcc 4.6.3 compiler with optimization flag set to $O3$. The Table I shows the hardware features of this machine. All implementations were written in C++. To compare the performance of different implementations, we considered the following parameters:

1. **Maximum Tree Size:** This depends on the size of the key space. We consider five different key ranges: 1000(1K), 10,000 (10K), 100,000 (100K), 1 million (1M) and 10 million (10M) keys. To capture only the steady state behaviour we *pre-populated* the tree to 50% of its maximum size, prior to starting the simulation run.
2. **Relative Distribution of Various Operations:** We consider three different workload distributions: (a) *read-dominated* workload : 90% search, 9% insert and 1% delete, (b) *mixed* workload : 70% search, 20% insert and 10% delete and (c) *write-dominated* workload : 0% search, 50% insert and 50% delete

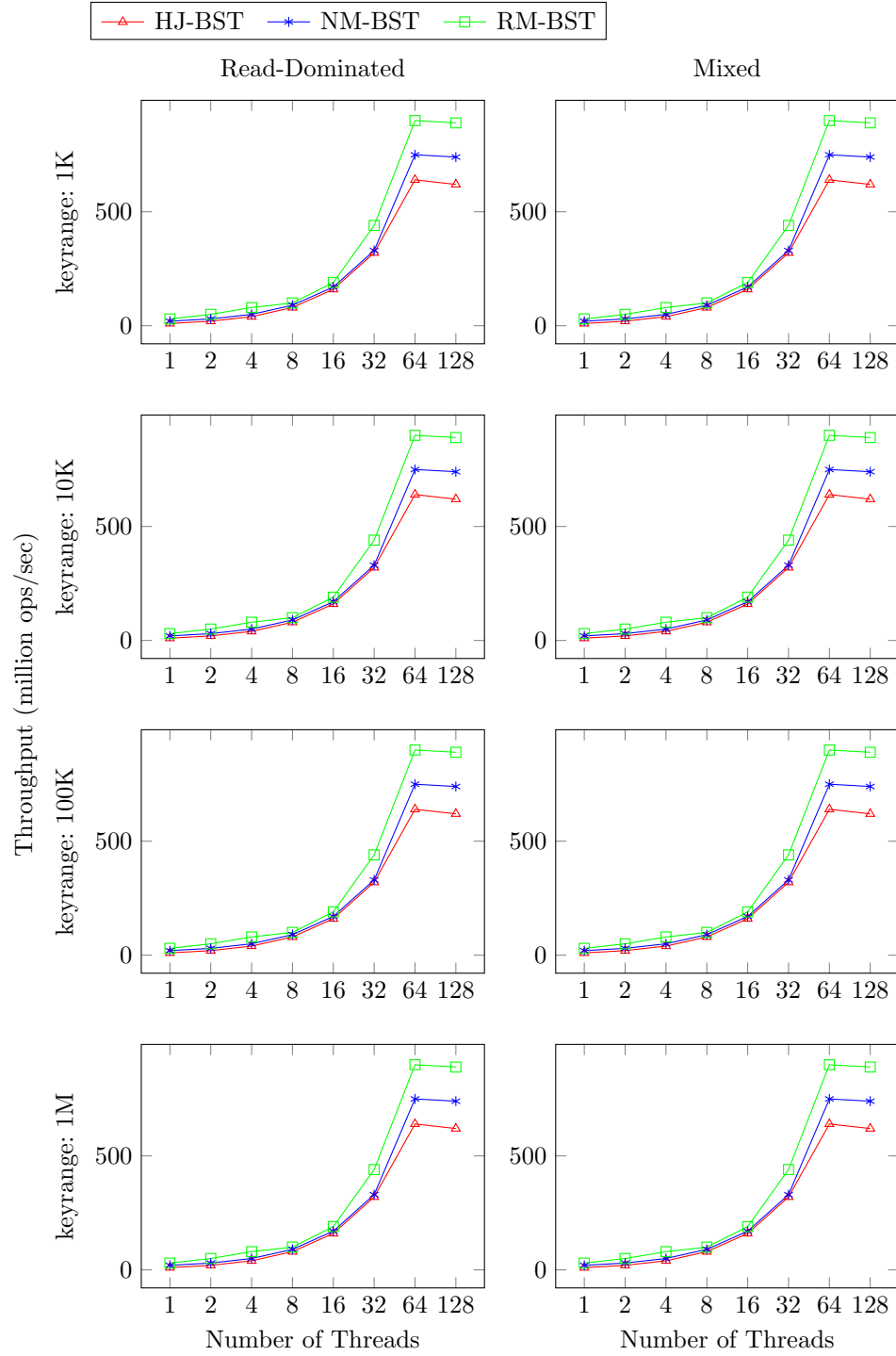
3. **Maximum Degree of Contention:** This depends on number of threads concurrently operating on the tree. We varied the number of threads from 1 to 128 in increments in powers of 2.

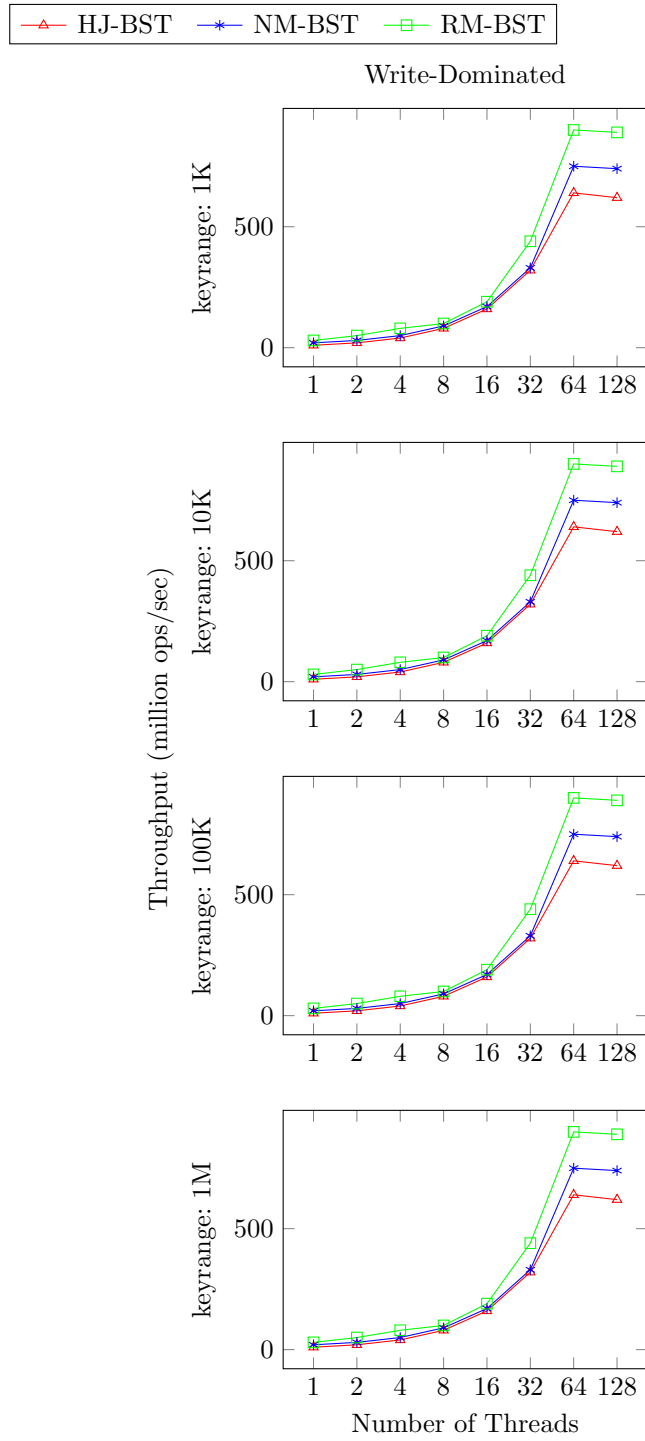
Table 1. hardware features

CPU packages	4	Algorithm		# of Objects Allocated		#of Atomic Instructions Executed	
Cores per package	8			Ins	Del	Ins	Del
Threads per core	2	Ellen & <i>et al.</i>		4	2	3	4
Clock frequency	2.3 GHz	Howley & Jones		2	1	3	upto 9
L1 cache (I/D)	64KB/16KB	Natarajan & Mittal		2	0	1	3
L2 cache	2 MB	This work		1	1	1	upto 6
L3 cache	6 MB						
Memory	256 GB						

We compared the performance of different implementations with respect to two metrics:

1. **System Throughput:** it is defined as the number of operations executed per unit time.
2. **Avg Seek Length:** it is defined as the average length of the access path of a seek operation.





References

1. Ellen, F., Fataourou, P., Ruppert, E., van Breugel, F.: Non-Blocking Binary Search Trees. In: Proceedings of the 29th ACM Symposium on Principles of Distributed Computing (PODC). pp. 131–140 (Jul 2010)
2. Howley, S.V., Jones, J.: A Non-Blocking Internal Binary Search Tree. In: Proceedings of the 24th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA). pp. 161–171 (Jun 2012)
3. Natarajan, A., Mittal, N.: Fast concurrent lock-free binary search trees. In: Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming. pp. 317–328. PPOPP '14, ACM, New York, NY, USA (2014), <http://doi.acm.org/10.1145/2555243.2555256>