

## A Speculation-Friendly Binary Search Tree

Tyler Crain, Vincent Gramoli, Michel Raynal

► **To cite this version:**

Tyler Crain, Vincent Gramoli, Michel Raynal. A Speculation-Friendly Binary Search Tree. [Research Report] PI-1984, 2011, pp.21. <inria-00618995v2>

**HAL Id: inria-00618995**

**<https://hal.inria.fr/inria-00618995v2>**

Submitted on 5 Mar 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## A Speculation-Friendly Binary Search Tree\*

Tyler Crain\*\*, Vincent Gramoli\*\*\*, Michel Raynal\*\*\*\*  
*tyler.crain@irisa.fr, vincent.gramoli@epfl.ch, raynal@irisa.fr*

**Abstract:** We introduce the first binary search tree algorithm designed for speculative executions. Prior to this work, tree structures were mainly designed for their pessimistic (non-speculative) accesses to have a bounded complexity. Researchers tried to evaluate transactional memory using such tree structures whose prominent example is the red-black tree library developed by Oracle Labs that is part of multiple benchmark distributions. Although well-engineered, such structures remain badly suited for speculative accesses, whose step complexity might raise dramatically with contention.

We show that our *speculation-friendly tree* outperforms the existing transaction-based version of the AVL and the red-black trees. Its key novelty stems from the *decoupling* of update operations: they are split into one transaction that modifies the abstraction state and multiple ones that restructure its tree implementation in the background. In particular, the speculation-friendly tree is shown correct, reusable and it speeds up a transaction-based travel reservation application by up to 3.5×.

**Key-words:** Transactional Memory, concurrent data structures, balanced binary trees

---

### *Un arbre binaire de recherche dédié aux accès transactionnels*

**Résumé :** *Les transactions, qui utilisent une technique de synchronisation optimiste, simplifient la programmation des architectures multi-cœurs. En effet, un programmeur a seulement besoin d'insérer des opérations dans des transactions afin d'obtenir un programme concurrent correct. Les programmeurs ont donc tout naturellement utilisé cette encapsulation sur plusieurs structures de données originellement dédiées à la programmation pessimiste (non optimistes) pour évaluer les transactions. L'exemple le plus probant étant l'arbre transactionnel rouge-noir développé par Oracle Labs et présent dans plusieurs distributions.*

*Malheureusement, ces structures de données ne sont pas adaptées à des exécutions optimistes car elles reposent sur des invariants contraignants, comme la profondeur logarithmique d'un arbre, pour borner le coût des accès pessimistes. Une telle complexité ne s'applique pas aux accès optimistes et garantir de tels invariants n'est pas nécessaire et induit de la contention en augmentant la probabilité pour une transaction d'avorter et de recommencer.*

*Dans ce papier, nous présentons un arbre binaire de recherche qui viole de façon transitoire ces invariants pour gagner en efficacité. Nous montrons que cet arbre est plus efficace que les arbres AVL et rouge-noir. Sa nouveauté majeure réside dans la séparation de ses opérations de modifications : elles sont coupées en une transaction qui modifie l'abstraction et plusieurs autres qui restructurent son implémentation. La bibliothèque qui en résulte limite la quantité d'avortements en augmentant légèrement le nombre d'accès, en particulier, elle améliore une application de réservation de voyage basée sur les transactions par un facteur multiplicatif de 3,5.*

**Mots clés :** *mémoire transactionnelle, arbre binaire, structures de données concurrente*

---

\* A 2-column 10 page version appeared at the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP 2012), ACM press.

\*\* Projet ASAP : équipe commune avec l'INRIA, le CNRS, l'université de Rennes 1

\*\*\* EPFL

\*\*\*\* Projet ASAP : équipe commune avec l'INRIA, le CNRS, l'université de Rennes 1; Senior member, Institut Universitaire de France

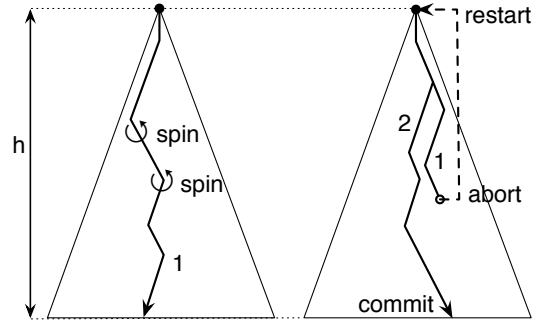


Figure 1: A balanced search tree whose complexity, in terms of the amount of accessed elements, is **(left)** proportional to  $h$  in a pessimistic execution and **(right)** proportional to the number of restarts in an optimistic execution

## 1 Introduction

The multicore era is changing the way we write concurrent programs. In such context, concurrent data structures are becoming a bottleneck building block of a wide variety of concurrent applications. Generally, they rely on invariants [31] which prevent them from scaling with multiple cores: a tree must typically remain sufficiently balanced at any point of the concurrent execution.

New programming constructs like transactions [18, 32] promise to exploit the concurrency inherent to multicore architectures. Most transactions build upon *optimistic synchronization*, where a sequence of shared accesses is executed speculatively and might abort. They simplify concurrent programming for two reasons. First, the programmer only needs to delimit regions of sequential code into transactions or to replace critical sections by transactions to obtain a safe concurrent program. Second, the resulting transactional program is reusable by any programmer, hence a programmer composing operations from a transactional library into another transaction is guaranteed to obtain new deadlock-free operations that execute atomically. By contrast, *pessimistic synchronization*, where each access to some location  $x$  blocks further accesses to  $x$ , is harder to program with [29, 30] and hampers reusability [15, 17].

Yet it is unclear how one can adapt a data structure to access it efficiently through transactions. As a drawback of the simplicity of using transactions, the existing transactional programs spanning from low level libraries to topmost applications directly derive from sequential or pessimistically synchronized programs. The impacts of optimistic synchronization on the execution is simply ignored.

To illustrate the difference between optimistic and pessimistic synchronizations consider the example of Figure 1 depicting their step complexity when traversing a tree of height  $h$  from its root to a leaf node. On the left, steps are executed pessimistically, potentially spinning before being able to acquire a lock, on the path converging towards the leaf node. On the right, steps are executed optimistically and some of them may abort and restart, depending on concurrent thread steps. The pessimistic execution of each thread is guaranteed to execute  $O(h)$  steps, yet the optimistic one may need to execute  $\Omega(hr)$  steps, where  $r$  is the number of restarts. Note that  $r$  depends on the probability of conflicts with concurrent transactions that depends, in turn, on the transaction length and  $h$ . Although it is clear that a transaction must be aborted before violating the abstraction implemented by this tree, e.g., inserting  $k$  successfully in a set where  $k$  already exists, it is unclear whether a transaction must be aborted before slightly unbalancing the tree implementation to strictly preserve the balance invariant.

We introduce a *speculation-friendly* tree as a tree that transiently breaks its balance structural invariant without hampering the abstraction consistency in order to speed up transaction-based accesses. Here are our contributions.

- We propose a speculation-friendly binary search tree data structure implementing an associative array and a set abstractions and decoupling the operations that modify the abstraction (we call these *abstract transactions*) from operations that modify the tree structure itself but not the abstraction (we call these *structural transactions*). An abstract transaction either inserts or deletes an element from the abstraction and in certain cases the insertion might also modify the tree structure. Some structural transactions rebalance the tree by executing a distributed rotation mechanism: each of these transactions executes a local rotation involving only a constant number of neighboring nodes. Some other structural transactions unlink and free a node that was logically deleted by a former abstract transaction.
- We prove the correctness (i.e., linearizability) of our tree and we compare its performance against existing transaction-based versions of an AVL tree and a red-black tree, widely used to evaluate transactions [7, 10, 12, 13, 19, 20, 33]. The speculation-friendly tree improves by up to  $1.6\times$  the performance of the AVL tree on the micro-benchmark and by up to  $3.5\times$  the performance of the built-in red-black tree on a travel reservation application, already well-engineered for transactions. Finally, our speculation-friendly tree performs similarly to a non-rotating tree but remains robust in face of non-uniform workloads.
- We illustrate (i) the portability of our speculation-friendly tree by evaluating it on two different Transactional Memories (TMs), TinySTM [13] and  $\mathcal{E}$ -STM [14] and with different configuration settings, hence outlining that our performance benefit is inde-

pendent from the transactional algorithm it uses; and (ii) its reusability by composing straightforwardly the remove and insert into a new move operation. In addition, we compare the benefit of relaxing data structures into speculation-friendly ones against the benefit of only relaxing transactions, by evaluating elastic transactions. It shows that, for a particular data structure, refactoring its algorithm is preferable to refactoring the underlying transaction algorithm.

The paper is organized as follows. In Section 2 we describe the problem related to the use of transactions in existing balanced trees. In Section 3 we present our speculation-friendly binary search tree. In Section 5 we evaluate our tree experimentally and illustrate its portability and reusability. In Section 6 we describe the related work and Section 7 concludes.

## 2 The Problem with Balanced Trees

In this section, we focus our attention on the structural invariant of existing tree libraries, namely the *balance*, and enlighten the impact of their restructuring, namely the *rebalancing*, on contention.

Trees provide logarithmic access time complexity given that they are balanced, meaning that among all downward paths from the root to a leaf, the length of the shortest path is not far apart the length of the longest path. Upon tree update, if their difference exceeds a given threshold, the structural invariant is broken and a rebalancing is triggered to restructure accordingly. This threshold depends on the considered algorithm: AVL trees [1] do not tolerate the longest length to exceed the shortest by 2 whereas red-black trees [4] tolerate the longest to be twice the shortest, thus restructuring less frequently. Yet in both cases the restructuring is triggered immediately when the threshold is reached to hide the imbalance from further operations.

Generally, one takes an existing tree algorithm and encapsulates all its accesses within transactions to obtain a concurrent tree whose accesses are guaranteed atomic (i.e., linearizable), however, the obtained concurrent transactions likely *conflict* (i.e., one accesses the same location another is modifying), resulting in the need to abort one of these transactions which leads to a significant waste of efforts. This is in part due to the fact that encapsulating an *update* operation (i.e., an insert or a remove operation) into a transaction boils down to encapsulating four phases in the same transaction:

1. the modification of the abstraction,
2. the corresponding structural adaptation,
3. a check to detect whether the threshold is reached and
4. the potential rebalancing.

**A transaction-based red-black tree** An example is the transaction-based binary search tree developed by Oracle Labs (formerly Sun Microsystems) and other researchers to extensively evaluate transactional memories [7, 10, 12, 13, 19, 20, 33]. This library relies on the classical red-black tree algorithm that bounds the step complexity of pessimistic insert/delete/contains. It has been slightly optimized for transactions by removing sentinel nodes to reduce false-conflicts, and we are aware of two benchmark-suite distributions that integrate it, STAMP [7] and synchrobench<sup>1</sup>.

Each of its update transactions encapsulate all the four phases given above even though phase (1) could be decoupled from phases (3) and (4) if transient violations of the balance invariant were tolerated. Such a decoupling is appealing given that phase (4) is subject to conflicts. In fact, the algorithm balances the tree by executing rotations starting from the position where a node is inserted or deleted and possibly going all the way up to the root. As depicted in Figure 2(a) and (b), a rotation consists of replacing the node where the rotation occurs by the child and adding this replaced node to one of its subtrees. A node cannot be accessed concurrently by an abstract transaction and a rotation, otherwise the abstract transaction might miss the node it targets while being rotated downward. Similarly, rotations cannot access common nodes as one rotation may unbalance the others.

Moreover, the red-black tree does not allow any abstract transaction to access a node that is concurrently being deleted from the abstraction because phases (1) and (2) are tightly coupled within the same transaction. If this was allowed the abstract transaction might end up on the node that is no longer part of the tree. Fortunately, if the modification operation is a deletion then phase (1) can be decoupled from the structural modification of phase (2) by marking the targeted node as logically deleted in phase (1) effectively removing it from the set abstraction prior to unlinking it physically in phase (2). This improvement is important as it lets a concurrent abstract transaction travel through the node concurrently being logically deleted in phase (1) without conflicting. Making things worse, without decoupling these four phases, having to abort within phase (4) would typically require the three previous phases to restart as well. Finally without decoupling only contains operations are guaranteed not to conflict with each other. With decoupling, insert/delete/contains do not conflict with each other unless they terminate on the same node as described in Section 3.

To conclude, for the transactions to preserve the atomicity and invariants of such a tree algorithm, they typically have to keep track of a large *read set* and *write set*, i.e., the sets of accessed memory locations that are protected by a transaction. Possessing large read/write sets increases the probability of conflicts and thus reduces concurrency. This is especially problematic in trees because the distribution

<sup>1</sup><http://lpd.epfl.ch/gramoli/php/synchrobench.php>

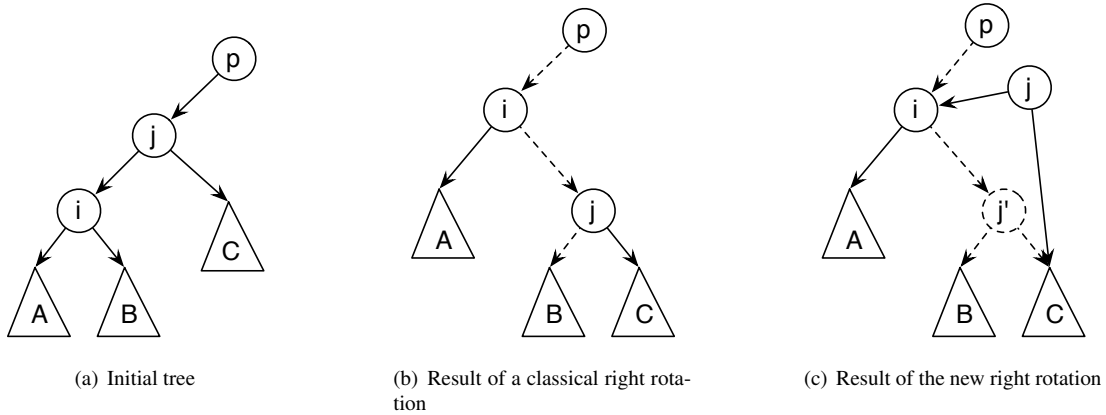


Figure 2: The classical rotation modifies node  $j$  in the tree and forces a concurrent traversal at this node to backtrack; the new rotation left  $j$  unmodified, adds  $j'$  and postpones the physical deletion of  $j$

of nodes in the read/write set is skewed so that the probability of the node being in the set is much higher for nodes near the root and the root is guaranteed to be in the read set.

**Illustration** To briefly illustrate the effect of tightly coupling update operations on the step complexity of classical transactional balanced trees we have counted the maximal number of reads necessary to complete typical insert/remove/contains operations. Note that this number includes the reads executed by the transaction each time it aborts in addition to the read set size of the transaction obtained at commit time.

Update	0%	10%	20%	30%	40%	50%
AVL tree	29	415	711	1008	1981	2081
Oracle red-black tree	31	573	965	1108	1484	1545
Speculation-friendly tree	29	75	123	120	144	180

Table 1: Maximum number of transactional reads per operation on three  $2^{12}$ -sized balanced search trees as the update ratio increases

We evaluated the aforementioned red-black tree, an AVL tree, and our speculation-friendly tree on a 48-core machine using the same transactional memory (TM) algorithm<sup>2</sup>. The expectation of the tree sizes is fixed to  $2^{12}$  during the experiments by performing an insert and a remove with the same probability. Table 1 depicts the maximum number of transactional reads per operation observed among 48 concurrent threads as we increase the update ratio, i.e., the proportion of insert/remove operations over contains operations.

For all three trees, the transactional read complexity of an operation increases with the update ratio due to the additional aborted efforts induced by the contention. Although the red-black and the AVL trees objective is to keep the complexity of pessimistic accesses  $O(\log_2 n)$  (proportional to 12 in this case), where  $n$  is the tree size, the read complexity of optimistic accesses grows significantly ( $14\times$  more at 10% update than at 0%, where there are no aborts) as the contention increases. As described in the sequel, the speculation-friendly tree succeeds in limiting the step complexity raise ( $2.6\times$  more at 10% update) of data structure accesses when compared against the transactional versions of state-of-the-art tree algorithms. An optimization further reducing the number of transactional reads between 2 (for 10% updates) and 18 (for 50% updates) is presented in Section 3.3.

### 3 The Speculation-Friendly Binary Search Tree

We introduce the speculation-friendly binary search tree by describing its implementation of an associative array abstraction, mapping a key to a value. In short, the tree speeds up the access transactions by decoupling two conflict-prone operations: the node deletion and the tree rotation. Although these two techniques have been used for decades in the context of data management [11, 24], our algorithm novelty lies in applying their combination to reduce transaction aborts. We first depict, in Algorithm 1, the pseudocode that looks like sequential code encapsulated within transactions before presenting, in Algorithm 2, more complex optimizations.

<sup>2</sup>TinySTM-CTL, i.e., with lazy acquirement [13].

**Algorithm 1** A Portable Speculation-Friendly Binary Search Tree

---

```

1: State of node  $n$ :
2:  $node$  a record with fields:
3:  $k \in \mathbb{N}$ , the node key
4:  $v \in \mathbb{N}$ , the node value
5:  $\ell, r \in \mathbb{N}$ , left/right child pointers, initially  $\perp$ 
6:  $left-h, right-h \in \mathbb{N}$ , local height of left/right
7: child, initially 0
8:  $local-h \in \mathbb{N}$ , expected local height, initially 1
9:  $del \in \{\text{true}, \text{false}\}$ , indicate whether
10: logically deleted, initially false

11: State of process  $p$ :
12:  $root$ , shared pointer to root

13:  $find(k)_p$ :
14:  $next \leftarrow root$ 
15: while true do
16:  $curr \leftarrow next$ 
17:  $val \leftarrow curr.k$ 
18: if  $val = k$  then break
19: if  $val > k$  then  $next \leftarrow read(curr.r)$ 
20: else  $next \leftarrow read(curr.l)$ 
21: if  $next = \perp$  then break
22: return  $curr$ 

23:  $contains(k)_p$ :
24: transaction {
25:  $result \leftarrow \text{true}$ 
26:  $curr \leftarrow find(k)$ 
27: if  $curr.k \neq k$  then  $result \leftarrow \text{false}$ 
28: else if  $read(curr.del)$  then  $result \leftarrow \text{false}$ 
29: } // current transaction tries to commit
30: return  $result$ 

31:  $insert(k, v)_p$ :
32: transaction {
33:  $result \leftarrow \text{true}$ 
34:  $curr \leftarrow find(k)$ 
35: if  $curr.k = k$  then
36: if  $read(curr.del)$  then  $write(curr.del, \text{false})$ 
37: else  $result \leftarrow \text{false}$ 
38: else // allocate a new node
39:  $new.k \leftarrow k$ 
40:  $new.v \leftarrow v$ 
41: if  $curr.k > k$  then  $write(curr.r, new)$ 
42: else  $write(curr.l, new)$ 
43: } // current transaction tries to commit
44: return  $result$ 

45:  $right\_rotate(parent, left-child)_p$ :
46: transaction {
47: if left-child then  $n \leftarrow read(parent.l)$ 
48: else  $n \leftarrow read(parent.r)$ 
49: if  $n = \perp$  then return false
50:  $\ell \leftarrow read(n.l)$ 
51: if  $\ell = \perp$  then return false
52:  $\ell r \leftarrow read(\ell.r)$ 
53:  $write(n.l, \ell r)$ 
54:  $write(\ell.r, n)$ 
55: if left-child then  $write(parent.l, \ell)$ 
56: else  $write(parent.r, \ell)$ 
57:  $update\_balance\_values()$ 
58: } // current transaction tries to commit
59: return true

60:  $delete(k)_p$ :
61: transaction {
62:  $result \leftarrow \text{true}$ 
63:  $curr \leftarrow find(k)$ 
64: if  $curr.k \neq k$  then
65:  $result \leftarrow \text{false}$ 
66: else
67: if  $read(curr.del)$  then  $result \leftarrow \text{false}$ 
68: else  $write(curr.del, \text{true})$ 
69: } // current transaction tries to commit
70: return  $result$ 

71:  $remove(parent, left-child)_p$ :
72: transaction {
73: if left-child then
74:  $n \leftarrow read(parent.l)$ 
75: else
76:  $n \leftarrow read(parent.r)$ 
77: if  $n = \perp$  or  $\neg read(n.del)$  then return false
78: if  $(child \leftarrow read(n.l)) \neq \perp$  then
79: if  $(child \leftarrow read(n.r)) \neq \perp$  then return false
80: if left-child then
81:  $write(parent.l, child)$ 
82: else
83:  $write(parent.r, child)$ 
84:  $update\_balance\_values()$ 
85: } // current transaction tries to commit
86: return true

```

---

### 3.1 Decoupling the tree rotation

The motivation for rotation decoupling stems from two separate observations: (i) a rotation is tied to the modification that triggers it, hence the process modifying the tree is also responsible for ensuring that its modification does not break the balance invariant and (ii) a rotation affects different parts of the tree, hence an isolated conflict can abort the rotation performed at multiple nodes. In response to these two issues we introduce a dedicated rotator thread to complete the modifying transactions faster and we distribute the rotation in multiple (node-)local transactions. Note that our rotating thread is similar to the collector thread proposed by Dijkstra et al. [11] to garbage collect stale nodes.

This decoupling allows the read set of the insert/delete operations to only contain the path from the root to the node(s) being modified and the write set to only contain the nodes that need to be modified in order to ensure the abstraction modification (i.e., the nodes at the bottom of the search path), thus reducing conflicts significantly. Let us consider a specific example. If rotations are performed within the insert/delete operations then each rotation increases the read and write set sizes. Take an insert operation that triggers a right rotation such as the one depicted in Figures 2(a)-2(b). Before the rotation the read set for the nodes  $p, j, i$  is  $\{p.l, j.r\}$ , where  $\ell$  and  $r$  represent the left and right pointers, and the write set is  $\emptyset$ . Now with the rotation the read set becomes  $\{p.l, i.r, j.l, j.r\}$  and the write set becomes  $\{p.l, i.r, j.l\}$  as denoted in the figure by dashed arrows. Due to  $p.l$  being modified, any concurrent transaction that traverses any part of this section of the tree (including all nodes  $i, j$ , and subtrees  $A, B, C, D$ ) will have a read/write conflict with this transaction. In the worst case an insert/delete operation triggers rotations all the way up to the root resulting in conflicts with all concurrent transactions.

**Rotation** As previously described, rotations are not required to ensure the atomicity of the insert/delete/contains operations so it is not necessary to perform rotations in the same transaction as the insert or delete. Instead we dedicate a separate thread that continuously checks for unbalances and rotates accordingly within its own node-local transactions.

More specifically, neither do the insert/delete operations comprise any rotation, nor do the rotations execute on a large block of nodes. Hence, local rotations that occur near the root can still cause a large amount of conflicts, but rotations performed further down the tree are less subject to conflict. If local rotations are performed in a single transaction block then even the rotations that occur further down the tree will be part of a likely conflicting transaction, so instead each local rotation is performed as a single transaction. Keeping the insert/delete/contains and rotate/remove operations as small as possible allows more operations to execute at the same time without conflicts, increasing concurrency.

Performing local rotations rather than global ones has other benefits. If rotations are performed as blocks then, due to concurrent insert/delete operations, not all of the rotations may still be valid once the transaction commits. Each concurrent insert/delete operation might require a certain set of rotations to balance the tree, but because the operations are happening concurrently the appropriate rotations to balance the tree are constantly changing and since each operation only has a partial view of the tree it might not know what the appropriate rotations are. With local rotations, each time a rotation is performed it uses the most up-to-date local information avoiding repeating rotations at the same location.

The actual code for the rotation is straightforward. Each rotation is performed just as it would be performed in a sequential binary tree (see Figure 2(a)-2(b)), but within a transaction.

Deciding when to perform a rotation is done based on local balance information omitted from the pseudocode. This technique was introduced in [5] and works as follows. *left-h* (resp. *right-h*) is a node-local variable to keep track of the estimated height of the left (resp. right) subtree. *local-h* (also a node-local variable) is always 1 larger than the maximum value of *left-h* and *right-h*. If the difference between *left-h* and *right-h* is greater than 1 then a rotation is triggered. After the rotation these values are updated as indicated by a dedicated function (line 57). Since these values are local to the node the estimated heights of the subtrees might not always be accurate. The propagate operation (described in the next paragraph) is used to update the estimated heights. Using the propagate operation and local rotations, the tree is guaranteed to be eventually perfectly balanced as in [5, 6].

**Propagation** The rotating thread executes continuously a depth-first traversal to propagate the balance information. Although it might propagate an outdated height information due to concurrency, the tree gets eventually balanced. The only requirement is that a node knows when it has an empty subtree (i.e., when *node.l* is  $\perp$ , *node.left-h* must be 0). This requirement is guaranteed since a new node is always added to the tree with *left-h* and *right-h* set to 0 and these values are updated when a node is removed or a rotation takes place. Each propagate operation is performed as a sequence of distributed transactions each acting on a single node. Such a transaction first travels to the left and right child nodes, checking their *local-h* values and using these values to update *left-h*, *right-h*, and *local-h* of the parent node. As no abstract transactions access these three values, they never conflict with propagate operations (unless the transactional memory used is inherently prone to false-sharing).

**Limitations** Unfortunately, spreading rotations and modifications into distinct transactions still does not allow insert/delete/contains operations that are being performed on separate keys to execute concurrently. Consider a delete operation that deletes a node at the root. In order to remove this node a successor is taken from the bottom of the tree so that it becomes the new root. This now creates a point of contention at the root and where the successor was removed. Every concurrent transaction that accesses the tree will have a read/write conflict with this transaction. Below we discuss how to address this issue.

### 3.2 Decoupling the node deletion

The speculation-friendly binary search tree exploits logical deletion to further reduce the amount of transaction conflicts. This two-phase deletion technique has been previously used for memory management like in [24], for example, to reduce locking in database indexes. Each node has a *deleted* flag, initialized to false when the node is inserted into the tree. First, the delete phase consists of removing the given key *k* from the abstraction—it logically deletes a node by setting a *deleted* flag to true (line 68). Second, the remove phase physically deletes the node from the tree to prevent it from growing too large. Each of these are performed as a separate transaction and the rotating thread is also responsible for garbage collecting nodes (cf. Section 3.4).

The deletion decoupling reduces conflicts by two means. First, it spreads out the two deletion phases in two separate transactions, hence reducing the size of the delete transaction. Second, deleting logically node *i* simply consists in setting the *deleted* flag to true (line 68), thus avoiding conflicts with concurrent abstract transactions that have traversed *i*.

**Find** The find operation is a helper function called implicitly by other functions within a transaction, thus it is never called explicitly by the application programmer. This operation looks for a given key *k* by parsing the tree similarly to a sequential code. At each node it goes right if the key of the node is larger than *k* (line 19), otherwise it goes left (line 20). Starting from the root it continues until it either finds a node with *k* (line 18) or until it reaches a leaf (line 21) returning the node (line 22). Notice that if it reaches a leaf, it has performed a transactional read on the child pointer of this leaf (lines 19–20), ensuring that some other concurrent transaction will not insert a node with key *k*.

**Contains** The contains operation first executes the find starting from the root, this returns a node (line 26). If the key of the node returned is equal to the key being searched for, then it performs a transactional read of the *deleted* flag (line 28). If the flag is false the operation returns true, otherwise it returns false. If the key of the returned node is not equal to the key being searched for then a node with the key being searched for is not in the tree and false is returned (lines 27 and 30).

**Insertion** The  $\text{insert}(k, v)$  operation uses the find procedure that returns a node (line 34). If a node is found with the same  $key$  as the one being searched for then the  $deleted$  flag is checked using a transactional read (line 36). If the flag is false then the tree already contains  $k$  and false is returned (lines 37 and 44). If the flag is true then the flag is updated to false (line 36) and true is returned. Otherwise if the  $key$  of the node returned is not equal to  $k$  then a new node is allocated and added as the appropriate child of the node returned by the find operation (lines 38-42). Notice that only in this final case does the operation modify the structure of the tree.

**Logical deletion** The delete uses also the find procedure in order to locate the node to be deleted (line 63). A transactional read is then performed on the  $deleted$  flag (line 67). If  $deleted$  is true then the operation returns false (lines 67 and 70), if  $deleted$  is false it is set to true (line 68) and the operation returns true. If the find procedure does not return a node with the same  $key$  as the one being searched for then false is returned (line 65 and 70). Notice that this operation never modifies the tree structure.

Consequently, the insert/delete/contains operations can only conflict with each other in two cases.

1. Two insert/delete/contains operations are being performed concurrently on some key  $k$  and a node with key  $k$  exists in the tree. Here (if at least one of the operations is an insert or delete) there will be a read/write conflict on the node's  $deleted$  flag. Note that there will be no conflict with any other concurrent operation that is being done on a different key.
2. An insert that is being performed for some key  $k$  where no node with key  $k$  exists in the tree. Here the insert operation will add a new node to the tree, and will have a read/write conflict with any operation that had read the pointer when it was  $\perp$  (before it was changed to point to the new node).

---

### Algorithm 2 Optimizations to the Speculation-Friendly Binary Search Tree

---

```

1: State of node  $n$ :
2:    $node$  the same record with an extra field:
3:    $rem \in \{\text{false}, \text{true}, \text{true\_by\_left\_rot}\}$ 
4:   indicate whether physically deleted,
5:   initially false

6:  $\text{remove}(parent, left-child)_p$ :
7:   transaction {
8:     if  $\text{read}(parent.rem)$  then return false
9:     if left-child then
10:       $n \leftarrow \text{read}(parent.l)$ 
11:     else
12:       $n \leftarrow \text{read}(parent.r)$ 
13:     if  $n = \perp$  or  $\neg \text{read}(n.deleted)$  then return false
14:     if  $(child \leftarrow \text{read}(n.l)) \neq \perp$  then
15:       if  $\text{read}(n.r) \neq \perp$  then return false
16:     else
17:        $child \leftarrow \text{read}(n.r)$ 
18:     if left-child then
19:        $\text{write}(parent.l, child)$ 
20:     else
21:        $\text{write}(parent.r, child)$ 
22:      $\text{write}(n.l, parent)$ 
23:      $\text{write}(n.r, parent)$ 
24:      $\text{write}(n.rem, \text{true})$ 
25:      $\text{update-balance-values}()$ 
26:   } // current transaction tries to commit
27:   return true

28:  $\text{find}(k)_p$ :
29:    $curr \leftarrow root$ 
30:    $next \leftarrow root$ 
31:   while true do
32:     while true do
33:        $rem \leftarrow \text{false}$ 
34:        $parent \leftarrow curr$ 
35:        $curr \leftarrow next$ 
36:        $val \leftarrow curr.k$ 
37:       if  $val = k$  then
38:         if  $\neg (rem \leftarrow \text{read}(curr.rem))$  then break
39:       if  $val > k \cup rem = \text{true\_by\_left\_rot}$  then
40:          $next \leftarrow \text{uread}(curr.r)$ 
41:       else  $next \leftarrow \text{uread}(curr.l)$ 
42:       if  $next = \perp$  then
43:         if  $\neg (rem \leftarrow \text{read}(curr.rem))$  then
44:           if  $val > k$  then  $next \leftarrow \text{read}(curr.r)$ 
45:           else  $next \leftarrow \text{read}(curr.l)$ 
46:         if  $next = \perp$  then break
47:       else
48:         if  $val \leq k$  then  $next \leftarrow \text{uread}(curr.r)$ 
49:         else  $next \leftarrow \text{uread}(curr.l)$ 
50:       if  $curr.k > parent.k$  then  $tmp \leftarrow \text{read}(parent.r)$ 
51:       else  $tmp \leftarrow \text{read}(parent.l)$ 
52:       if  $curr = tmp$  then
53:         break
54:       else
55:          $next \leftarrow curr$ 
56:          $curr \leftarrow parent$ 
57:       return curr

58:  $\text{right\_rotate}(parent, left-child)_p$ :
59:   transaction {
60:     if  $\text{read}(parent.rem)$  then
61:       return false
62:     if left-child then
63:        $n \leftarrow \text{read}(parent.l)$ 
64:     else
65:        $n \leftarrow \text{read}(parent.r)$ 
66:     if  $n = \perp$  then
67:       return false
68:      $l \leftarrow \text{read}(n.l)$ 
69:     if  $l = \perp$  then
70:       return false
71:      $lr \leftarrow \text{read}(l.r)$ 
72:      $r \leftarrow \text{read}(n.r)$ 
73:     // allocate a new node
74:      $new.k \leftarrow n.k$ 
75:      $new.l \leftarrow lr$ 
76:      $new.r \leftarrow r$ 
77:      $\text{write}(l.r, new)$ 
78:      $\text{write}(n.rem, \text{true})$ 
79:     // In the case of a left rotate set
80:     //  $n.rem$  to  $\text{true\_by\_left\_rot}$ 
81:     if left-child then
82:        $\text{write}(parent.l, l)$ 
83:     else
84:        $\text{write}(parent.r, l)$ 
85:      $\text{update-balance-values}()$ 
86:   } // current transaction tries to commit
87:   return true

```

---

**Physical removal** Removing a node that has no children is as simple as unlinking the node from its parent (lines 81–83). Removing a node that has 1 child is done by just unlinking it from its parent, then linking its parent to its child (also lines 81–83). Each of these removal procedures is a very small transaction, only performing a single transactional write. This transaction conflicts only with concurrent transactions that read the link from the parent before it is changed.

Upon removal of a node  $i$  with two children, the node in the tree with the immediately larger key than  $i$ 's must be found at the bottom of the tree. This performs reads on all the way to the leaf and a write at the parent of  $i$ , creating a conflict with any operation that has traversed this node. Fortunately, in practice such removals are not necessary. In fact only nodes with no more than one child are



removed from the tree (if the node has two children, the remove operation returns without doing anything, cf. line 79). It turns out that removing nodes with no more than one children is enough to keep the tree from growing so large that it affects performance.

The removal operation is performed by the maintenance thread. While it is traversing the tree performing rotation and propagate operations it also checks for logically deleted nodes to be removed.

**Limitations** The traversal phase of most functions is prone to false-conflicts, as it comprises read operations that do not actually need to return values from the same snapshot. Specifically, by the time a traversal transaction reaches a leaf, the value it read at the root likely no longer matters, thus a conflict with a concurrent root update could simply be ignored. Nevertheless, the standard TM interface forces all transactions to adopt the same strongest semantics prone to false-conflicts [15]. In the next paragraphs we discuss how to extend the basic TM interface to cope with such false-conflicts.

### 3.3 Optional improvements

In previous sections, we have described a speculation-friendly tree that fulfills the standard TM interface [21] for the sake of portability across a large body of research work on TM. Now, we propose to further reduce aborts related to the rotation and the find operation at the cost of an additional lightweight read operation, `uread`, that breaks this interface. This optimization is thus usable only in TM systems providing additional explicit calls and do not aim at replacing but complementing the previous algorithm to preserve its portability. This optimization complementing Algorithm 1 is depicted in Algorithm 2, it does not affect the existing `contains/insert/delete` operations besides speeding up their internal find operation. Here the left rotation is not the exact symmetry of the right rotation code.

This change to the algorithm requires that each node has an additional flag indicating whether or not the node has been physically removed from the tree (a node is physically removed during a successful rotate or remove operation). This removed flag can be set to false, true or `true_by_left_rot` and is initialized to false. In order not to complicate the pseudo code `true_by_left_rot` is considered to be equivalent to true, only on line 39 of the find operation is this parameter value specifically checked for.

**Lightweight reads** The key idea is to avoid validating superfluous read accesses when an operation traverses the tree structure. This idea has been exploited by elastic transactions that use a bounded buffer instead of a read set to validate only immediately preceding reads, thus implementing a form of hand-over-hand locking transaction for search structure [14]. We could have used different extensions to implement these optimizations. DSTM [20] proposes early release to force a transaction stop keeping track of a read set entry. Alternatively, the current distribution of TinySTM [13] comprises unit loads that do not record anything in the read set. While we could have used any of these approaches to increase concurrency we have chosen the unit loads of TinySTM, hence the name `uread`. This `uread` returns the most recent value written to memory by a committed transaction by potentially spin-waiting on the location until it stops being concurrently modified.

A first interesting result, is that the read/write set sizes can be kept at a size of  $O(k)$  instead of the  $O(k \log n)$  obtained with the previous tree algorithm, where  $k$  is the number of nested `contains/insert/delete` operations nested in a transaction. The reasoning behind this is as follow: Upon success, a `contains` only needs to ensure that the node it found is still in the tree when the transaction commits, and can ignore the state of other nodes it had traversed. Upon failure, it only needs to ensure that the node  $i$  it is looking for is not in the tree when the transaction commits, this requires to check whether the pointer from the parent that would point to  $i$  is  $\perp$  (i.e., this pointer should be in the read set of the transaction and its value is  $\perp$ ). In a similar vein, `insert` and `delete` only need to validate the position in the tree where they aimed at inserting or deleting. Therefore, `contains/insert/delete` only increases the size of the read/write set by a constant instead of a logarithmic amount.

It is worth mentioning that `ureads` have a further advantage over normal reads other than making conflicts less likely: Classical reads are more expensive to perform than unit reads. This is because in addition to needing to store a list keeping track of the reads done so far, an opaque TM that uses invisible reads needs to perform validation of the read set with a worst case cost of  $O(s^2)$ , where  $s$  is the size of the read set, whereas a TM that uses visible reads performs a modification to shared memory for each read.

**Rotation** Rotations remain conflict-prone in Algorithm 1 as they incur a conflict when crossing the region of the tree traversed by a `contains/insert/delete` operation. If `ureads` are used in the `contains/insert/delete` operations then rotations will only conflict with these operations if they finish at one of the two nodes that are rotated by rotation operation (for example in Figure 2(a) this would be the node  $i$  or  $j$ ). A rotation at the root will only conflict with a `contains/insert/delete` that finished at (or at the rotated child of) the root, any operations that travel further down the tree will not conflict.

Figure 2(c) displays the result of the new rotation that is slightly different than the previous one. Instead of modifying  $j$  directly,  $j$  is unlinked from its parent (effectively removing it from the tree, lines 82–84) and a new node  $j'$  is created (line 73), taking  $j$ 's place in the tree (lines 82–84). During the rotation  $j$  has a removed flag that is set to true (line 78), letting concurrent operations know that  $j$  is no longer in the tree but its deallocation is postponed. Now consider a concurrent operation that is traversing the tree and is preempted on  $j$  during the rotation. If a normal rotation is performed the concurrent operation will either have to backtrack or the transaction would have to abort (as the node it is searching for might be in the subtree  $A$ ). Using the new rotation, the preempted operation will still have a path to  $A$ .

As previous noted the removed flag can be set to one of three values (false, true or true\_by\_left\_rot). Only when a node is removed during a left rotation is the flag set to true\_by\_left\_rot. This is necessary to ensure that the find operation follows the correct path in the specific case that the operation is preempted on a node that is concurrently removed by a left rotation and this node has the same key  $k$  as the one being searched for. In this case the find operation must travel to the right child of the removed node otherwise it might miss the node with key  $k$  that has replaced the removed node from the rotation. In all other cases the find operation can follow the child pointer as normal.

**Find, contains and delete** The interesting point for the find operation is that the search continues until it finds a node with the *removed* flag set to false (line 38 and 43). Once the leaf or a node with the same key as the one being searched for is reached, a transactional read is performed on the *removed* flag to ensure that the node is not removed from the tree (by some other operation) at least until the transaction commits. If *removed* is true then the operation continues traversing the tree, otherwise the correct node has been found. Next, if the node is a leaf, a transactional read must be performed on the appropriate child pointer to ensure this node remains a leaf throughout the transaction (lines 44–45). If this read does not return  $\perp$  then the operation continues traversing the tree. Otherwise the operation then leaves the nested while loop (lines 38 and 46), but the find operation does not return yet.

One additional transactional read must be performed to ensure safety. This is the read of the parent's pointer to the node about to be returned (lines 50–51). If this read does not return the same node as found previously, the find operation continues parsing the tree starting from the parent (lines 55–56). Otherwise the process leaves the while loop (line 53) and the node is returned (line 57). By performing a transactional read on the parent's pointer we ensure the STM system performs a validation before continuing.

The advantage of this updated find operation is that ureads are used to traverse the tree, it only uses transactional reads to ensure atomicity when it reaches what is suspected to be the last node it has to traverse. The original algorithm exclusively uses transactional reads to traverse the tree and because of this, modifications to the structure of the tree that occur along the traversed path cause conflicts, which do not occur in the updated algorithm. The contains/insert/delete operations themselves are identical in both algorithms.

**Removal** The remove operation requires some modification to ensure safety when using ureads during the traversal phase. Normally if a contains/insert/delete operation is preempted on a node that is removed then that operation will have to backtrack or abort the transaction. This can be avoided as follows. When a node is removed, its left and right child pointers are set to point to its previous parent (lines 22–23). This provides a preempted operation with a path back to the tree. The removed node also has its *removed* flag set to true (line 24) letting preempted operations know it is no longer in the tree (the node is left to be freed later by garbage collection).

### 3.4 Garbage collection

As explained previously, there is always a single rotator thread that continuously executes a recursive depth first traversal. It updates the local, left and right heights of each node and performs a rotation or removal if necessary. Nodes that are successfully removed are then added to a garbage collection list. Each application thread maintains a boolean indicating a pending operation and a counter indicating the number of completed operations. Before starting a traversal, the rotator thread sets a pointer to what is currently the end of the garbage collection list and copies all booleans and counters. After a traversal, if for every thread its counter has increased or if its boolean is false then the nodes up to the previously stored end pointer can be safely freed. Experimentally, we found that the size of the list was never larger than a small fraction of the size of the tree but theoretically we expect the total space required to remain linear in the tree size.

## 4 Correctness Proof

There are two parts in the proof. First we have to ensure the structure of the tree is always a valid binary search tree. This is important because in a binary search tree at any time there is exactly one correct location for a key  $k$ , the term used for such a tree is *valid binary tree*. A binary tree that does not have the previous property is simply called a *binary tree*. Second we have to show the insert, delete, and contains operations are linearizable.

It is important to remember for this proof that when a transaction commits the transactional reads and writes appear as if they have happened atomically and that unit-read operations only return values from previously committed transactions (or the value written by the current transaction, if the transaction has written to the location being read).

Another important part of this proof is the way the tree is first created. It is created with a root node with key  $\infty$  so that all nodes will always be on its left subtree. This node will always be the root (i.e. it will not be modified by rotate or removal operations), this makes simpler operations and proofs.

**Binary trees** Each node has two boolean state variables. When the variable *deleted* is false it entails that the value of *node.key* is in the set represented by the tree. When it is true it entails that the value of *node.key* is not in the set represented by the tree. When the variable *removed* is false it entails that the node exists in the tree (meaning a path exists from the root of the tree to the node). When it is

true (or `true_by_left_rot`) it entails that the node does not exist in the tree (meaning no path exists from the root of the tree to the node). For simplicity throughout the pseudo code `true_by_left_rot` is considered to be equivalent to `true`, only on line 39 of the `find` operation is this parameter value specifically checked for.

In the proof we will use the phrase *a node can reach a range of keys* which is explained here. Take the root, from this node there is a path to any node in the tree meaning any key that is in the set is reachable from the root. Furthermore for any key that is not in the tree the root has a path to where it would be inserted (the leaf that would be its parent). This means that the root with key  $k$  node has a range  $[-\infty, \infty]$ . Now its left child has range  $[-\infty, k]$  and its right child has range  $[k, -\infty]$ . Or for example consider some node with range  $(10, 20]$  and key 14. Its left child will have a range  $(10, 14)$  and its right child will have a range  $(14, 20]$ . have a path to it.

The phrase *a node  $n$  has a path to a range of keys at least as large as some other node  $n'$*  is also used in the proof. It means that every key that is in the range of  $n'$  is also in the range of  $n$  (and possibly some more). For example any node will have a path to a range of keys at least as large as its left child (in fact it has the exact range of its left and right child combined).

**Set operations** Here traditional operations on the set are defined in the context of transactions. It is important to remember that the TM guarantees a linearization of transactions.

The following definitions are used in defining the operations. Saying a key  $k$  is in (not in) the set before the committal of a transaction  $T$  means that if some transaction  $T1$  performs a contains operation on key  $k$  and is serialized as the transaction immediately before  $T$ ,  $T1$  would return true (false).

Saying a key  $k$  is in (not in) the set after the committal of a transaction  $T$  means that if some transaction  $T2$  performs a contains operation on key  $k$  and is serialized as the transaction immediately after  $T$ ,  $T2$  would return true (false).

**delete** For transaction that commits a successful delete operation of key  $k$ , before the commit  $k$  was in the set and afterwards  $k$  is not in the set. For transaction that commits a failed delete operation of key  $k$ , before and after the commit  $k$  was not in the set.

**insert** For transaction that commits a successful insert operation of key  $k$ , before the commit  $k$  was not in the set and after the commit  $k$  is in the set. For transaction that commits a failed insert operation of key  $k$ , before and after the commit  $k$  is in the set.

**contains** For transaction that commits a successful contains operation of key  $k$ , before and after the commit  $k$  was in the set. For transaction that commits a failed contains operation of key  $k$ , before and after the commit  $k$  was not in the set.

**Lemma 1** *A node has at most one parent with `removed = false`.*

**Proof.** Assume by contradiction that a node has more than one parent with `removed = false`. There are three operation where a node can be given a new parent. First during the insert operations on lines 41–42, but this is a new node so before this line it has no parent. Second during the remove operation on lines 19–21, but by line 8 the other parent has `removed` set to true. Third during the *right-rotate* operation the node  $l$  gets a new parent on lines 82–84, but by line 78 the other parent has `removed` set to true. Also during the *right-rotate* operation the node  $n2$  gets  $l$  as a parent (line 77), but since it is a new node it has no other parent. (This holds for the *left-rotate* operation by symmetry) Given this, a node will have at most one parent with `removed = false`.  $\square$

**Lemma 2** *A node with `removed = false` only has paths from it to other nodes with `removed = false`.*

**Proof.** This proof is by induction on the number  $j$  of operations done on a node with key  $k$ .

The base case is when a new node is created and added to the tree. This can happen in the insert operation. During the operation a new node is created with no children and for itself it has `removed = false` (line 38) and the proof holds.

Now for the induction step from  $j = m - 1$  to  $j = m$ , this could be a contains, delete, insert, remove, or rotate operation. First note that the contains and delete operations do not change any children pointers. If this is an insert operation then at  $j = m - 1$  *left* or *right* must be  $\perp$  (line 46 of the `find` operation) and the proof obviously still holds. If this is a remove operation, then the child of the node is being removed. This means that the new child will either be  $\perp$  or the child of the child (lines 14–21), but by induction these nodes have `removed = false`. By symmetry this holds for the right and left children. Otherwise this is a rotate operation. First consider right rotations. By induction we know that node  $n$  only has paths to nodes with `removed = false`. After the rotation node  $l$  points to nodes that had paths from  $n$  before the rotation as well as  $n2$  (line 77) which has `removed = false`.  $n2$  points to nodes that had paths from  $n$  before the rotation (lines 75–76). By symmetry this holds for left rotations.  $\square$

**Lemma 3** *A node with `removed = false` has at least one parent with `removed = false` that points to it (except the root, as it has no parent).*

**Proof.** Assume by contradiction that there is no parent with `removed = false`. When a node is first added to the tree it has a parent with `removed = false` (line 43 of the `find` operation). The only operations that removes links from nodes are the remove (line 24) and rotate (line 78) operations, but in each case when a link is removed, a new link from a different node with `removed = false` is added (lines 19–21 of remove and 82–84 of rotate).  $\square$

**Lemma 4** *A node with `removed = false` has a path from the root node to it.*

**Proof.** Given that the root is always the same node and it always has  $removed = false$  the proof of this lemma follows directly from lemmas 2 and 3.  $\square$

**Lemma 5** *A node with  $removed = true$  has no path from the root node to it.*

**Proof.** By the way the tree is structured the root node always has  $removed = false$ . Now it follows directly from lemma 2 that there is no path from the root to a node with  $removed = true$ .  $\square$

**Lemma 6** *The nodes with  $removed = false$  make up a single binary tree.*

**Proof.** From the structure of a node, it can have at most 2 children. Now by lemmas 1, 4, and 5 the proof follows.  $\square$

**Lemma 7** *A rotation operation on a valid binary search tree results in a valid binary tree of the nodes with  $removed = false$ .*

**Proof.** From Figure 2 which describes the *rotation* operation and due to the use of transactional reads/writes we can see that the resulting tree is equivalent to a tree with a classical binary tree rotation performed on it.  $\square$

For the following Lemma, we assume that the find operation is performed correctly (this will be proved in a later lemma).

**Lemma 8** *Assuming a correct find operation, a successful insert operation on a valid binary search tree results in a valid binary tree of the nodes with  $removed = false$ .*

**Proof.** There are two cases to consider.

First the key  $k$  that we are inserting is already contained in the tree with  $removed = false$  (lines 37–38 of find). In this case the structure of the tree will not be modified, thus the resulting tree will still be valid.

Second consider that there is no node in the tree with key  $k$ . In this case the find operation will return the correct node that will then become the parent of the new node with key  $k$  (line 41–42). Using transactional reads, the find operation ensures that the parent node has  $removed = false$  (line 43) and  $\perp$  (line 46) for the child pointer where the new node will be inserted. The new node is then created and added to the tree as the child of the node returned from find. This is done using a transactional write (lines 41–42) which ensures that the value of the pointer was  $\perp$  before the write, and finally resulting in a valid tree containing the new node after the transaction commits.  $\square$

**Lemma 9** *A successful remove operation on a valid binary search tree results in a valid binary tree that does not contain the node removed.*

**Proof.** A removal can only be performed on a node  $n$  with at least one  $\perp$  child which is ensured by a transactional read (lines 14–17). Node  $n$  being removed is unlinked from its parent (lines 19–21) (the parent is ensured to be in the tree by a transaction read on  $removed = false$ ) effectively removing it from the tree (lemma 1). If both children of  $n$  are set to  $\perp$ , then the parent's new child becomes  $\perp$  (lines 14–21) leaving the tree still valid. If node  $n$  has a child  $c$  such that  $c \neq \perp$ , then  $c$  becomes the parent's new child (lines 14–21). By lemma 2 this  $c$  must have  $removed = false$  and by lemma 4 it is part of the valid binary tree during the transaction (until the transaction commits). Thus the resulting tree is still valid.  $\square$

**Lemma 10** *Modifications to the tree structure are only performed on nodes with  $removed = false$ .*

**Proof.** A rotate, insert or remove operation can modify the tree.

During a right rotate operation the nodes that are modified are  $n$ ,  $l$  and the parent of  $n$  (lines 77–84). A transactional read is performed on the remove variable of the parent node (line 61), this along with lemma 2 ensures that  $removed = false$  for the parent of  $n$  as well as  $n$ , and  $l$ . By symmetry the left rotate operation also only modifies nodes with  $removed = false$ .

During a successful insert operation a new node might be added to the tree, in this case its parent node is modified (lines 41–42), and a transactional read is performed on the parent to ensure that  $removed = false$  (line 43 of the find operation).

During a remove operation a node is removed from the tree. A transactional read is performed on the node and its parent (line 8), this along with lemma 2 ensures that  $removed = false$  for both nodes.  $\square$

**Lemma 11** *From a node with  $removed = true$  there is always a path to some node with  $removed = false$ .*

**Proof.** There are two places where a node can be set to  $removed = true$ . First during the remove operation, in this case the node that is removed has both of the nodes child pointers are set to a node with  $removed = false$  (lines 22–23). Second during the rotate operation, in this case the node that is removed does not have its child pointers changed. By line 70  $n$  must have at least 1 child and using lemma 2 this child must have  $removed = false$ .

Now notice that once the  $removed$  field of a node is set to true it will never be reverted to false (lemma 10). This along with the use of induction on the length of the path to a node with  $removed = false$  completes the proof.  $\square$

**Lemma 12** *From any node every path leads to a leaf node (or  $\perp$ ).*

**Proof.** This proof is the same as lemma 11 with a small modification.

First consider a node with  $removed = false$ . By lemma 6 it is clear that there is a path from this node to a leaf node.

Now consider a node with  $removed = true$ . There are two places where a node's  $removed$  flag can be set to true. It can either be set in the remove operation, but in this case the node's left and right pointers are set to a node with  $removed = false$  (lines 22–23).

Or it can be set in the rotation operation, first notice that in this case the node's left and right pointers are not changed. Then before the rotation the node has  $removed = false$  and therefore by lemma 2 the (node's left and right) pointers will point to either nodes with  $removed = false$  or  $\perp$ .

Now that after a node has had  $removed$  set to true the node will not be modified again (lemma 10), this along with lemma 6 and the use of induction on the length of the path to a node with  $removed = false$  or  $\perp$  completes the proof.  $\square$

The phrase *a node that has  $removed = true$  from a rotation operation* refers to the node that is no longer in the tree after the rotation. For example in figure 2(c) this would be the node  $j$ . This phrase is used in the following lemmas.

**Lemma 13** *A node  $n$  with key  $k$  that has  $removed = true$  (true\_by\_left\_rot) from a right\_rotation operation for its **left** child has a path to a range of keys at least as large as  $n$  did before the rotation (including the new node  $n2$  with key  $k$ ), and for its **right** child a path to a range of keys at least as large as the right child before the rotation, or  $\perp$ .*

**Proof.** Assume that the node  $n$  has key  $k$  and before the rotation has a path to a range  $[a, b]$ . This means that node  $l$  (the left child of  $n$ ) has a path to the range  $[a, k]$ . Assume node  $l$  has key  $j$ . The right child of  $n$  is either  $bot$  or some node  $r$  with a path to the range  $[k, b]$ .

After the rotation  $l$  keeps its left child with its right child changing to  $n2$ .  $n2$ 's right child becomes  $n$ 's right child, and  $n2$ 's left child becomes  $l$ 's old right child. This leaves  $n2$  with a path to the range  $[j, b]$ , and  $l$  with a path to the range  $[a, b]$ .

During the *rotation* operation no modifications are made to node  $n$ , except setting  $removed = true$ . Now  $n$ 's left child is  $l$  giving it a path to the range of keys  $[a, b]$ . Node  $n$  still has the same right child who still has the same range,  $[k, b]$ .  $\square$

**Lemma 14** *A node  $n$  with key  $k$  that has  $removed = true$  from a left\_rotation operation for its **right** child has a path to a range of keys at least as large as  $n$  did before the rotation (including the new node  $n2$  with key  $k$ ), and for its **left** child a path to a range of keys at least as large as the left child before the rotation, or  $\perp$ .*

**Proof.** This proof follows from symmetry (the left rotation is the mirror of the right rotation) and lemma 13.  $\square$

**Lemma 15** *A node that has  $removed = true$  from a remove operation has a path to a range of keys as least as large as just before the remove operation took place.*

**Proof.** Assume that the node  $n$  (where  $n$  is the node to be removed) has key  $k$ . Assume that  $n$  has a parent node  $p$  with range  $[a, b]$  with key  $j$ . This leaves  $n$  with range  $[a, k]$  if  $n$  is the left child (or  $[k, b]$  if  $n$  is the right child, the proof of this case will follow by symmetry).

After the removal  $removed$  will be set to true (line 24 of remove) for  $n$  and both its child pointers will point to  $p$  (lines 22-23 of remove). Node  $p$  will still have a range of  $[a, b]$  and will be reachable from  $n$ , which completes the proof.  $\square$

**Lemma 16** *A node that has  $removed = true$  has either:*

- for each child a path to a range of keys at least as large as they did when it had  $removed = false$  or
- a single  $\perp$  child with the other child having a path to a range of keys at least as large as both children before when it had  $removed = false$ .

**Proof.** The proof follows using lemmas 13, 14, and 15 as well as induction on the length of the path from the node to a node with  $removed = false$ .  $\square$

**Lemma 17** *A find operation on a valid binary search tree always returns the correct location from the valid binary tree.*

**Proof.** Assume by contradiction that a find operation does not return the correct location from the valid binary tree. By lemma 6 we know that every node with  $removed = false$  is a node in the valid binary tree so there are two possibilities. The operation never returns or the operation returns the wrong location. First for the operation never returning. From lemma 12 we know that there are no cycles in the path from any node so the operation will not get stuck in an endless loop. In order for the operation to complete it must reach a node with  $removed = false$  that either matches the key being searched for (lines 37–38) or has  $\perp$  as the appropriate child (lines 43–45). Lemma 11 is enough to show that this will always happen and the operation will terminate.

The second possibility where the operation returns the wrong location is more difficult to prove. To this end, we prove by induction on the number of nodes traversed by the operation that the find operation will always have a path to the correct location. In order for the operation to reach the correct location the following must be true: After each traversal from one node to the next the operation must either be at the correct location or have a path from the current location to a range of keys that includes the key being searched for. Once the correct location is reached transactional reads are used to ensure that the location remains correct throughout the transaction (lines 38, and 43–45).

- The base case is easy given that the operation starts from the root which is always part of the valid tree and can reach all nodes with  $removed = false$  (lemma 4).
- Now the induction step. Assume that after  $n - 1$  nodes have been traversed the operation has a path to a range of keys that includes the key  $k$  that is being searched for. If the  $n - 1^{th}$  node is the correct location the the proof is done, otherwise the operation must travel to the  $n^{th}$  node. Now it must be shown that the after the operation travels to the  $n^{th}$  node it is still on the correct path. There are 2 possibilities.
  1. The operation can move from a node with  $removed = false$  (at the time of the load of the child pointer, lines 40–41). By lemma 2 the child must also have  $removed = false$  (at the time of the load of the child pointer), in this case the traversal is performed on a valid binary tree (lemma 6). Since the choice of the the child is based on a standard tree traversal the operation must still be on the correct path.
  2. The operation can move from a node with  $removed = true$ . By the assumption given by induction the  $n - 1^{th}$  node has a path to the range of keys that includes  $k$ . From the  $n - 1^{th}$  node either the right or left child must be chosen. The node is chosen based on a standard tree traversal with one exception: If the node has the same key  $k$  that is being searched for and the node was removed by a left rotation then the right child is chosen (lines 39–40). With this exception then in all cases if the node is not  $\perp$  then by lemma 16 it must have a path to the same range as when it has  $removed = false$ , which includes  $k$  (the  $n - 1^{th}$  node's range includes  $k$  so the  $n^{th}$  node's range must also). Otherwise if one child is  $\perp$  then the other child is chosen (lines 48–49), which must have a path to a range at least as large as the  $n - 1^{th}$  node by lemma 16 (which includes  $k$ ).

□

**Theorem 1** *An insert operation is valid.*

**Proof.** The insert operation starts by performing a find operation (line 34). From lemma 17 we know that the find operation will return the correct place. There are two possibilities, either the find operation returns a node in the tree with key  $k$  (line 37), or it returns a node in the tree that has  $\perp$  for the child where  $k$  must be inserted (line 46).

First consider the case where a node with key  $k$  is returned. The transactional read on  $removed$  (line 38 of the find operation) ensures that the node will be in the tree throughout the duration of the transaction. Next a transactional read is done on the node's  $deleted$  variable (line 36) ensuring that this value will remain the same throughout the transaction. If the read on  $deleted$  returns false then the insert operation can return false because the node exists in the set. Otherwise if the read on  $deleted$  returns true then a transaction performs a transactional write setting  $deleted$  to false (line 36). The transactional read on  $deleted$  ensures that  $k$  is not in the set before the transaction commits. The transactional write on  $deleted$  ensures that  $k$  is in the set after the transaction commits.

Second consider the case where a node with  $key \neq k$  is returned. On lines 44–45 of the find operation a transactional read has been done on the child pointer of this node, returning  $\perp$ , which must be valid throughout the transaction. The transactional read on  $removed$  (line 43 of the find operation) ensures that the node will be in the tree throughout the duration of the transaction. By lemma 6 this node belongs to a correct binary tree, this along with lemma 17 ensures that the child of this node is the only place where a node with key  $k$  could exist (and since the value of the child pointer of  $\perp$  a node with key  $k$  is not in the tree). A new node with key  $k$  is created and then added to the tree using a transactional write to the child pointer (lines 41–42). The  $removed$  and  $deleted$  fields of a newly created node can only be false so when the transaction commits the new node will be in the tree and  $k$  will be in the set. □

**Theorem 2** *A contains operation is valid.*

**Proof.** The contains operation starts by performing a find operation (line 26). From lemma 17 we know that the find operation will return the correct place. There are two possibilities, either the find operation returns a node in the tree with key  $k$  (line 37), or it returns a node in the tree that has  $\perp$  for the child where  $k$  could exist (line 43).

First consider the case where a node with key  $k$  is returned. The transactional read on *removed* (line 38 of the find operation) ensures that the node will be in the tree throughout the duration of the transaction. Next a transactional read is done on the node's *deleted* variable (line 28) ensuring that this value will remain the same throughout the transaction. If the read on *deleted* returns false then the *contains* operation can return true because the node exists in the set otherwise false can be returned because the node does not exist in the set.

Second consider the case where a node with  $key \neq k$  is returned. On lines 44–45 of the find operation a transactional read has been done on the child pointer of this node, returning  $\perp$ , and due to the transactional read the pointer must remain as  $\perp$  throughout the transaction. The transactional read on *removed* (line 43 of the find operation) ensures that the node will be in the tree throughout the duration of the transaction. By lemma 6 this node belongs to a correct binary tree, this along with lemma 17 ensures that the child of this node is the only place where a node with key  $k$  could exist (and since the value of the child pointer is  $\perp$  a node with key  $k$  is not in the tree). The *contains* operation can then return false.  $\square$

**Theorem 3** *A delete operation is valid.*

**Proof.** The delete operation is almost the same as the *contains* operation. The only difference is in the case where a node with key  $k$  is returned from the find operation. The transactional read on *removed* (line 38 of the find operation) ensures that the node will be in the tree throughout the duration of the transaction. Next a transactional read is done on the node's *deleted* variable (line 67) ensuring that this value will remain the same throughout the transaction. If the read on *deleted* returns true then the delete operation can return false because the node does not exist in the set. Otherwise  $k$  is in the set and a transactional write is performed setting the nodes *deleted* variable to true. Since *removed* is false this node is part of the valid tree so it is the only place where key  $k$  can be so by setting *deleted* to true  $k$  must not be in the set.  $\square$

## 5 Experimental Evaluation

We experimented our library by integrating it in (i) a micro-benchmark of the synchrobench suite to get a precise understanding of the performance causes and in (ii) the tree-based vacation reservation system of the STAMP suite and whose runs sometimes exceed half an hour. The machine used for our experiments is a four AMD Opteron 12-core Processor 6172 at 2.1 Ghz with 32 GB of RAM running Linux 2.6.32, thus comprising 48 cores in total.

### 5.1 Testbed choices

We evaluate our tree against well-engineered tree algorithms especially dedicated to transactional workloads. The red-black tree is a mature implementation developed and improved by expert programmers from Oracle Labs and others to show good performance of TM in numerous papers [7, 10, 12, 13, 19, 20, 33]. The observed performance is generally scalable when contention is low, most of integer set benchmarks on which they are tested consider the ratio of attempted updates instead of effective updates. To avoid the misleading (attempted) update ratios that capture the number of calls to potentially updating operations, we consider the *effective* update ratios of synchrobench counting only modifications and ignoring the operations that fail (e.g., remove may fail in finding its parameter value thus failing in modifying the data structure).

The AVL tree we evaluate (as well as the aforementioned red-black tree) is part of STAMP [7]. As mentioned before one of the main refactoring of this red-black tree implementation is to avoid the use of sentinel nodes that would produce false-conflicts within transactions. This improvement could be considered a first-step towards obtaining a speculation-friendly binary search tree, however, the modification-restructuring, which remains tightly coupled, prevents scalability to high levels of parallelism.

To evaluate performance we ran the micro-benchmark and the vacation application with 1, 4, 8, 16, 24, 32, 40, 48 application threads. For the micro-benchmark, we averaged the data over three runs of 10 seconds each. For the vacation application, we averaged the data over three runs as well but we used the recommended default settings and some runs exceeded half an hour because of the amount of transactions used. We carefully verified that the variance was sufficiently low for the result to be meaningful.

### 5.2 Biased workloads and the effect of restructuring

In this section, we evaluate the performance of our speculation-friendly tree on an integer set micro-benchmark providing remove, insert, and contains operations, similarly to the benchmark used to evaluate state-of-the-art TM algorithms [9, 13, 14]. We implemented two set libraries that we added to the synchrobench distribution: our non-optimized speculation-friendly tree and a baseline tree that is similar but never rebalances the structure whatever modifications occur. Figure 3 depicts the performance obtained from four different binary search trees: the red-black tree (RBtree), our speculation-friendly tree without optimizations (SFtree), the no-restructuring tree (NRtree) and the AVL tree (AVLtree).

The performance is expressed as the number of operations executed per microsecond. The update ratio varies between 5% and 20%. As we obtained similar results with  $2^{10}$ ,  $2^{12}$  and  $2^{14}$  elements, we only report the results obtained from an initialized set of  $2^{12}$  elements.

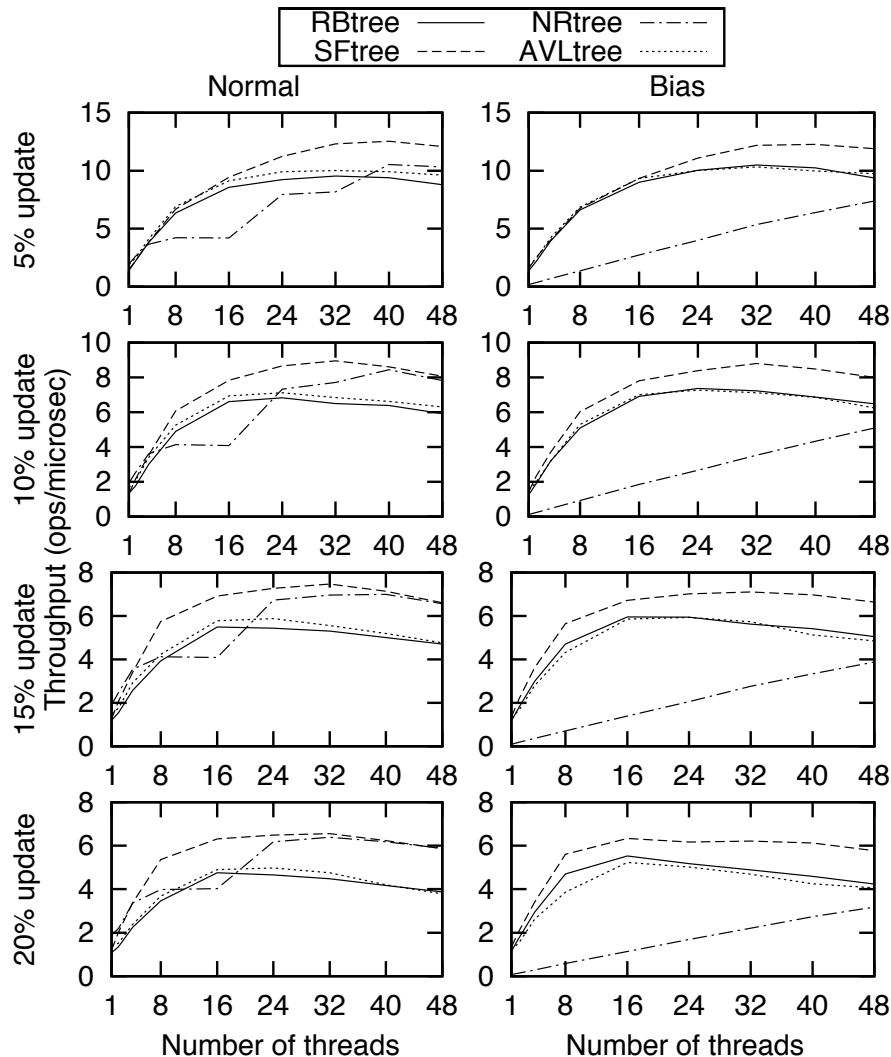


Figure 3: Comparing the AVL tree (AVLtree), the red-black tree (RBtree), the no-restructuring tree (NRtree) against the speculation-friendly tree (SFtree) on an integer set micro-benchmark with from 5% (**top**) to 20% updates (**bottom**) under normal (**left**) and biased (**right**) workloads

The biased workload consists of inserting (resp. deleting) random values skewed towards high (resp. low) numbers in the value range: the values always taken from a range of  $2^{14}$  are skewed with a fixed probability by incrementing (resp. decrementing) with an integer uniformly taken within  $[0..9]$ .

On both the normal (uniformly distributed) and biased workloads, the speculation-friendly tree scales well up to 32/40 threads. The no-restructuring tree performance drops to a linear shape under the biased workload as expected: as it does not rebalance, the complexity increases with the length of the longest path from the root to a leaf that, in turn, increases with the number of performed updates. In contrast, the speculation-friendly tree can only be unbalanced during a transient period of time which is too short to affect the performance even under biased workloads.

The speculation-friendly tree improves both the red-black tree and the AVL tree performance by up to  $1.5\times$  and  $1.6\times$ , respectively. The speculation-friendly tree is less prone to contention than AVL and red-black trees, which both share similar performance penalties due to contention.

### 5.3 Portability to other TM algorithms

The speculation-friendly tree is an inherently efficient data structure that is portable to any TM systems. It fulfills the TM interface standardized in [21] and thus does not require the presence of explicit escape mechanisms like early release [20] or snap [8] to avoid extra TM bookkeeping (our uread optimization being optional). Nor does it require high-level conflict detection, like open nesting [2, 25, 26]



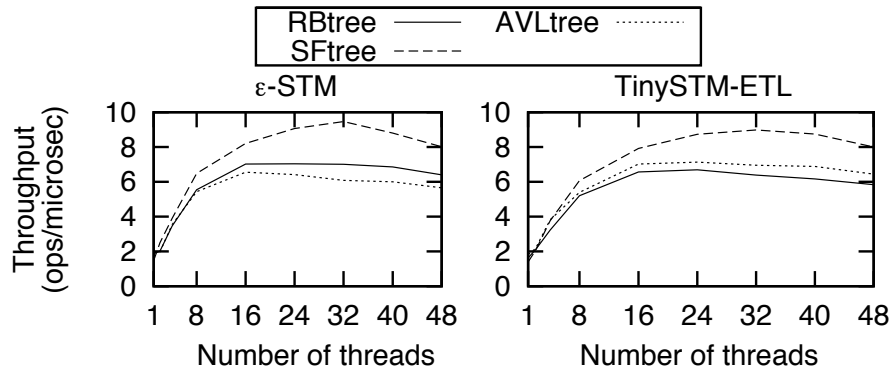


Figure 4: The speculation-friendly library running with **(left)** another TM library ( $\mathcal{E}$ -STM) and with **(right)** the previous TM library in a different configuration (TinySTM-ETL, i.e., with eager acquirement)

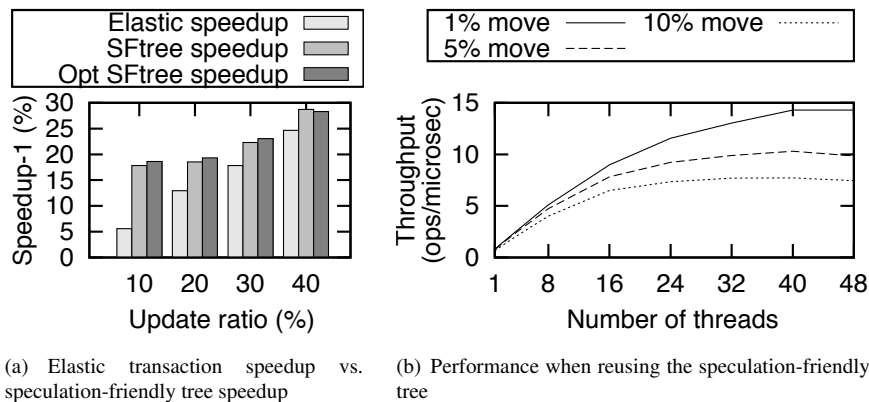


Figure 5: Elastic transaction comparison and reusability

or transactional boosting [19]. Such improvements rely on explicit calls or user-defined abstract locks, and are not supported by existing TM compilers [21] which limits their portability. To make sure that the obtained results are not biased by the underlying TM algorithm, we evaluated the trees on top of  $\mathcal{E}$ -STM [14], another TM library (on a  $2^{16}$  sized tree where  $\mathcal{E}$ -STM proved efficient), and on top of a different TM design from the one used so far: with eager acquirement.

The obtained results, depicted in Figure 4 look similar to the ones obtained with TinySTM-CTL (Figure 3) in that the speculation-friendly tree executes faster than other trees for all TM settings. This suggests that the improvement of speculation-friendly tree is potentially independent from the TM system used. A more detailed comparison of the improvement obtained using elastic transactions on red-black trees against the improvement of replacing the red-black tree by the speculation-friendly tree is depicted in Figure 5(a). It shows that the elastic improvements (15% on average) is lower than the speculation-friendly tree one (22% on average, be it optimized or not).

#### 5.4 Reusability for specific application needs

We illustrate the reusability of the speculation-friendly tree by composing remove and insert from the existing interface to obtain a new atomic and deadlock-free move operation. Reusability is appealing to simplify concurrent programming by making it modular: a programmer can reuse a library without having to understand its synchronization internals. While reusability of sequential programs is straightforward, concurrent programs can generally be reused only if the programmer understands how each element is protected. For example, reusing a library can lead to deadlocks if shared data are locked in a different order than what is recommended by the library. Additionally, a lock-striping library may not conflict with a concurrent program that locks locations independently even though they protect common locations, thus leading to inconsistencies.

Figure 5(b) indicates the performance on workloads comprising 90% of read-only operations (including contains and failed updates) and 10% move/insert/delete effective update operations (among which from 1% to 10% are move operations). The performance de-

creases as more move operations execute, because a move protects more elements in the data structure than a simple insert or delete operation and during a longer period of time.

## 5.5 The vacation travel reservation application

We experiment our optimized library tree with a travel reservation application from the STAMP suite [7], called *vacation*. This application is suitable for evaluating concurrent binary search tree as it represents a database with four tables implemented as tree-based directories (cars, rooms, flights, and customers) accessed concurrently by client transactions.

Figure 6 depicts the execution time of the STAMP *vacation* application building on the Oracle red-black tree library (by default), our optimized speculation-friendly tree, and the baseline no-restructuring tree. We added the speedup obtained with each of these tree libraries over the performance of bare sequential code of *vacation* without synchronization. (A concurrent tree library outperforms the sequential tree when its speedup exceeds 1.) The chosen workloads are the two default configurations (“low contention” and “high contention”) taken from the STAMP release, with the default number of transactions,  $8\times$  more transactions than by default and  $16\times$  more, to increase the duration and the contention of the benchmark without using more threads than cores.

*Vacation* executes always faster on top of our speculation-friendly tree than on top of its built-in Oracle red-black tree. For example, the speculation-friendly tree improves performance by up to  $1.3\times$  with the default number of transactions and to  $3.5\times$  with  $16\times$  more transactions. The reason of this is twofold: (i) In contrast with the speculation-friendly tree, if an operation on the red-black tree traverses a location that is being deleted then this operation and the deletion conflict. (ii) Even though the Oracle red-black tree tolerates that the longest path from the root to a leaf can be twice as long as the shortest one, it triggers the rotation immediately after this threshold is reached. By contrast, our speculation-friendly tree keeps checking the unbalance to potentially rotate in the background. In particular, we observed on 8 threads in the high contention settings that the red-black tree *vacation* triggered around 130,000 rotations whereas the speculation-friendly *vacation* triggered only 50,000 rotations.

Finally, we observe that *vacation* presents similarly good performance on top of the no-restructuring tree library. In rare cases, the speculation-friendly tree outperforms the no-restructuring tree probably because the no-restructuring tree does not physically remove nodes from the tree, thus leading to a larger tree than the abstraction. Overall, their performance is comparable. With  $16\times$  the default number of transactions, the contention gets higher and rotations are more costly.

## 6 Related Work

Aside from the optimistic synchronization context, various relaxed balanced trees have been proposed. The idea of decoupling the update and the rebalancing was originally proposed by Guibas and Sedgwick [16] and was applied to AVL trees by Kessels [22], and Nurmi, Soisalon-Soininen and Wood [28], and to red-black trees by Nurmi and Soisalon-Soininen [27]. Manber and Ladner propose a lock-based tree whose rebalancing is the task of separate maintenance threads running with a low priority [23]. Bougé et al. [5] propose to lock a constant number of nodes within local rotations. The combination of local rotations executed by different threads self-stabilizes to a tree where no nodes are marked for removal. The main objective of these techniques is still to keep the tree depth low enough for the lock-based operations to be efficient. Such solutions do not apply to speculative operations due to aborts.

Ballard [3] proposes a relaxed red-black tree insertion well-suited for transactions. When an insertion unbalances the red-black tree it marks the inserted node rather than rebalancing the tree immediately. Another transaction encountering the marked node must rebalance the tree before restarting. The relaxed insertion was shown generally more efficient than the original insertion when run with DSTM [20] on 4 cores. Even though the solution limits the waste of effort per aborting rotation, it increases the number of restarts per rotation. By contrast, our local rotation does not require the rotating transaction to restart, hence benefiting both insertions and removals.

Bronson et al. [6] introduce an efficient object-oriented binary search tree. The algorithm uses underlying time-based TM principles to achieve good performance, however, its operations cannot be encapsulated within transactions. For example, a key optimization of this tree distinguishes whether a modification at some node  $i$  grows or shrinks the subtree rooted in  $i$ . A conflict involving a growth could be ignored as no descendant are removed and a search preempted at node  $i$  will safely resume in the resulting subtree. Such an optimization is not possible using TMs that track conflicts between read/write accesses to the shared memory. This implementation choice results in higher performance by avoiding the TM overhead, but limits reusability due to the lack of bookkeeping. For example, a programmer willing to implement a size operation would need to explicitly clone the data structure to disable the growth optimization. Therefore, the programmer of a concurrent application that builds upon this binary search tree library must be aware of the synchronization internals of this library (including the growth optimization) to reuse it.

Felber, Gramoli and Guerraoui [14] specify the elastic transactional model that ignores false conflicts but guarantees reusability. In the companion technical report, the red-black tree library from Oracle Labs was shown executing efficiently on top of an implementation of the elastic transaction model,  $\mathcal{E}$ -STM. The implementation idea consists of encapsulating the (i) operations that locate a position in the red-black tree (like insert, contains, delete) into an *elastic* transaction to increase concurrency and (ii) other operations, like size, into a regular transaction. This approach is orthogonal to ours as it aims at improving the performance of the underlying TM, independently from the data structure, by introducing relaxed transactions. Hence, although elastic transactions can cut themselves upon conflict

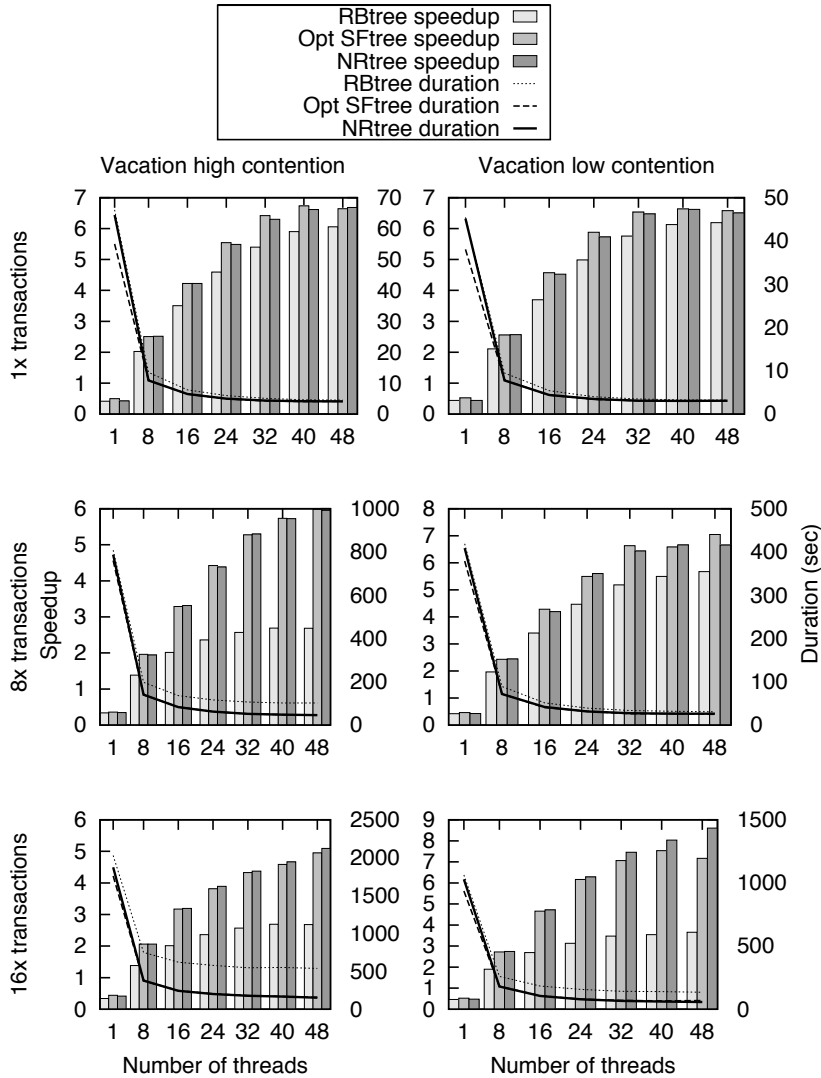


Figure 6: The speedup (over single-threaded sequential) and the corresponding duration of the vacation application built upon the red-black tree (RBtree), the optimized speculation-friendly tree (Opt SFtree) and the no-restructuring tree (NRtree) on **(left)** high contention and **(right)** low contention workloads, and with **(top)** the default number of transaction, **(middle)**  $8\times$  more transactions and **(bottom)**  $16\times$  more transactions

detection, the resulting  $\mathcal{E}$ -STM, still suffers from congestion and wasted work when applied to non-speculation-friendly data structures. The results presented in Section 5.3 confirm that the elastic speedup is even higher when the tree is speculation-friendly.

## 7 Conclusion

Transaction-based data structures are becoming a bottleneck in multicore programming, playing the role of synchronization toolboxes a programmer can rely on to write a concurrent application easily. This work is the first to show that speculative executions require the design of new data structures. The underlying challenge is to decrease the inherent contention by relaxing the invariants of the structure while preserving the invariants of the abstraction.

In contrast with the traditional pessimistic synchronization, the optimistic synchronization allows the programmer to directly observe the impact of contention as part of the step complexity because conflicts potentially lead to subsequent speculative re-executions. We

have illustrated, using a binary search tree, how one can exploit this information to design a speculation-friendly data structure. The next challenge is to adapt this technique to a large body of data structures to derive a speculation-friendly library.

## Source Code

The code of the speculation-friendly binary search tree is available at <http://lpd.epfl.ch/gramoli/php/synchrobench.php>.

## Acknowledgements

The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement number 238639, ITN project TransForm, and grant agreement number 248465, the S(o)OS project.

## References

- [1] G. Adelson-Velskii and E. M. Landis. An algorithm for the organization of information. In *Proc. of the USSR Academy of Sciences*, volume 146, pages 263–266, 1962.
- [2] Kunal Agrawal, I-Ting Angelina Lee, and Jim Sukha. Safe open-nested transactions through ownership. In *Proc. of the 14th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2009.
- [3] Lucia Ballard. Conflict avoidance: Data structures in transactional memory, May 2006. Undergraduate thesis, Brown University.
- [4] Rudolf Bayer. Symmetric binary b-trees: Data structure and maintenance algorithms. *Acta Informatica* 1, 1(4):290–306, 1972.
- [5] Luc Bougé, Joaquim Gabarro, Xavier Messeguer, and Nicolas Schabanel. Height-relaxed AVL rebalancing: A unified, fine-grained approach to concurrent dictionaries, 1998. Research Report 1998-18, ENS Lyon.
- [6] Nathan G. Bronson, Jared Casper, Hassan Chafi, and Kunle Olukotun. A practical concurrent binary search tree. In *Proc. of the 15th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2010.
- [7] Chi Cao Minh, JaeWoong Chung, Christos Kozyrakis, and Kunle Olukotun. STAMP: Stanford transactional applications for multi-processing. In *Proc. of The IEEE Int'l Symp. on Workload Characterization*, 2008.
- [8] Christopher Cole and Maurice Herlihy. Snapshots and software transactional memory. *Sci. Comput. Program.*, 58(3):310–324, 2005.
- [9] Luke Dalessandro, Michael Spear, and Michael L. Scott. NOrec: streamlining STM by abolishing ownership records. In *Proc. of the 15th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2010.
- [10] D. Dice, O. Shalev, , and N. Shavit. Transactional locking II. In *Proc. of the 20th Int'l Symp. on Distributed Computing*, 2006.
- [11] Edsger W. Dijkstra, Leslie Lamport, A. J. Martin, C. S. Scholten, and E. F. M. Steffens. On-the-fly garbage collection: an exercise in cooperation. *Commun. ACM*, 21(11):966–975, 1978.
- [12] Aleksandar Dragojevic, Pascal Felber, Vincent Gramoli, and Rachid Guerraoui. Why STM can be more than a research toy. *Commun. ACM*, 54(4):70–77, 2011.
- [13] Pascal Felber, Christof Fetzer, and Torvald Riegel. Dynamic performance tuning of word-based software transactional memory. In *Proc. of the 13th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2008.
- [14] Pascal Felber, Vincent Gramoli, and Rachid Guerraoui. Elastic transactions. In *Proc. of the 23rd Int'l Symp. on Distributed Computing*, 2009.
- [15] Vincent Gramoli and Rachid Guerraoui. Democratizing transactional programming. In *Proc. of the ACM/IFIP/USENIX 12th Int'l Middleware Conference*, 2011.
- [16] L. J. Guibas and R. Sedgwick. A dichromatic framework for balanced trees. In *Proc. of the 19th Annual Symp. on Foundations of Computer Science*, 1978.
- [17] Tim Harris, Simon Marlow, Simon Peyton-Jones, and Maurice Herlihy. Composable memory transactions. In *Proc. of the 10th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2005.

- [18] M. Herlihy and J. E. B. Moss. Transactional memory: Architectural support for lock-free data structures. In *Proc. of the 20th Annual Int'l Symp. on Computer Architecture*, 1993.
- [19] Maurice Herlihy and Eric Koskinen. Transactional boosting: A methodology for highly-concurrent transactional objects. In *Proc. of the 13th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2008.
- [20] Maurice Herlihy, Victor Luchangco, Mark Moir, and William N. Scherer, III. Software transactional memory for dynamic-sized data structures. In *Proc. of the 22nd Annual ACM SIGACT-SIGOPS Symp. on Principles of Distributed Computing*, 2003.
- [21] Intel Corporation. Intel transactional memory compiler and runtime application binary interface, May 2009.
- [22] J. L. W. Kessels. On-the-fly optimization of data structures. *Comm. ACM*, 26:895–901, 1983.
- [23] Udi Manbar and Richard E. Ladner. Concurrency control in a dynamic search structure. *ACM Trans. Database Syst.*, 9(3):439–455, 1984.
- [24] C. Mohan. Commit-LSN: a novel and simple method for reducing locking and latching in transaction processing systems. In *Proc. of the 16th Int'l Conference on Very Large Data Bases*, 1990.
- [25] J. Eliot B. Moss. Open nested transactions: Semantics and support. In *Workshop on Memory Performance Issues*, 2006.
- [26] Yang Ni, Vijay Menon, Ali-Reza Abd-Tabatabai, Antony L. Hosking, Richard L. Hudson, J. Eliot B. Moss, Bratin Saha, and Tatiana Shpeisman. Open nesting in software transactional memory. In *Proc. of the 12th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2007.
- [27] O. Nurmi and E. Soisalon-Soininen. Uncoupling updating and rebalancing in chromatic binary search trees. In *Proc. of the 10th ACM Symp. on Principles of Database Systems*, 1991.
- [28] O. Nurmi, E. Soisalon-Soininen, and D. Wood. Concurrency control in database structures with relaxed balance. In *Proc. of the 6th ACM Symp. on Principles of Database Systems*, 1987.
- [29] Victor Pankratius and Ali-Reza Adl-Tabatabai. A study of transactional memory vs. locks in practice. In *Proc. of the 23rd ACM Symp. on Parallelism in Algorithms and Architectures*, 2011.
- [30] Christopher J. Rossbach, Owen S. Hofmann, and Emmett Witchel. Is transactional programming actually easier? In *Proc. of the 15th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming*, 2010.
- [31] Nir Shavit. Data structures in the multicore age. *Commun. ACM*, 54(3):76–84, 2011.
- [32] Nir Shavit and Dan Touitou. Software transactional memory. In *Proc. of the 14th ACM Symp. on Principles of Distributed Computing*, 1995.
- [33] Richard M. Yoo, Yang Ni, Adam Welc, Bratin Saha, Ali-Reza Adl-Tabatabai, and Hsien-Hsin S. Lee. Kicking the tires of software transactional memory: why the going gets tough. In *Proc. of the 20th ACM Symp. on Parallelism in Algorithms and Architectures*, 2008.