# Analysis and Evaluation of Non-Blocking Interpolation Search Trees

Aleksandar Prokopec
Oracle Labs
Switzerland
aleksandar.prokopec@gmail.com

Trevor Brown
University of Waterloo
Canada
me@tbrown.pro

Dan Alistarh
Institute of Science and Technology
Austria
dan.alistarh@ist.ac.at

## Abstract

We start by summarizing the recently proposed implementation of the first non-blocking concurrent interpolation search tree (C-IST) data structure. We then analyze the individual operations of the C-IST, and show that they are correct and linearizable. We furthermore show that lookup (and several other non-destructive operations) are wait-free, and that the insert and delete operations are lock-free.

We continue by showing that the C-IST has the following properties. For arbitrary key distributions, this data structure ensures worst-case $O(\log n + p)$ amortized time for search, insertion and deletion traversals. When the input key distributions are *smooth*, lookups run in expected $O(\log \log n + p)$ time, and insertion and deletion run in expected amortized $O(\log \log n + p)$ time, where $p$ is a bound on the number of threads.

Finally, we present an extended experimental evaluation of the non-blocking IST performance.

***CCS Concepts*** • **Theory of computation** → **Concurrent algorithms**; *Shared memory algorithms*; • **Computing methodologies** → **Concurrent algorithms**;

***Keywords*** concurrent data structures, search trees, interpolation, non-blocking algorithms

## 1 Introduction

The recently proposed non-blocking interpolation search tree (C-IST) improves the asymptotic complexity of the basic tree operations, such as insertion, lookup and predecessor queries, from logarithmic to doubly-logarithmic for a specific class of workloads whose input key distributions are *smooth* [1]. For arbitrary key distributions, C-IST operations retain standard logarithmic running time bounds.

---

[1] Intuitively speaking, this property bounds the height of the "spikes" in the probability density function of the input key set. The work on sequential interpolation search trees defines smoothness more precisely [34].

While C-IST was proposed and described in detail in related work [15], this paper analyzes the C-IST data structure from the correctness and an asymptotic complexity standpoint. This work furthermore extends the experimental analysis of the runtime behavior and performance of a concrete implementation of C-IST, and compares C-IST to a wider range of data structures.

After summarizing the C-IST data structure and its operations in Section 2, we put forth the following contributions:

- We prove that the basic C-IST operations are correct, linearizable and non-blocking (wait-free lookup, and lock-free modifications) in Section 3.
- We analyze the asymptotic complexity of the C-IST algorithm in Section 4. In particular, we show that the worst-case cost of a lookup operation is $O(p + \log n)$, where $p$ is the bound on the number of concurrent threads and $n$ is the number of keys in the C-IST; and we show that the worst-case amortized running time of insert and delete operations is $O(\gamma(p + \log n))$, where $\gamma$ is a bound on average interval contention. Furthermore, we show that, when the input key distribution is smooth, the expected worst-case running time of a lookup operation is $O(p + \log \log n)$; and the expected amortized running time of insertion and deletion is $O(\gamma(p + \log \log n))$.
- We present an extended experimental evaluation of our C-IST implementation in Section 5. The evaluation includes a performance comparison against seven state-of-the-art concurrent tree data structures, a measurement of last-level cache misses between C-IST and best-performing competitors, a breakdown of execution time across different key-set sizes, parallelism levels and update rates, the impact of using a multi-counter in the root node, and a performance evaluation on a No-SQL database workload, and a Zipfian distribution workload, at different skew levels.

The paper ends with an overview of related work in Section 6, and a short conclusion in Section 7.

## 2 Non-Blocking Interpolation Search Tree

### 2.1 Data Types

Non-blocking IST consists of the data types shown in Figure 1. The main data type, named `IST`, represents the interpolation search tree. It has only one field called `root`, which is a pointer to the root node of the data structure. When the data structure is created, the `root` field points to an empty leaf node of type `Empty`. The data type `Single` is used as a leaf node that has a single key and an associated value. The data type `Inner` represents inner nodes in the tree. The relationship between these node types is shown on the right side of Figure 1.

The `Inner` data type contains the search keys, and the pointers to the children nodes. In addition to that, the `Inner` data type has the field `degree`, which holds the number of children of that inner node, and a field called `initSize`, which holds the number of keys the corresponding subtree had at the moment when this inner node was created. Apart from the pointers in the `children` array, all of these fields are written once when the inner node is created, and are not modified after that.

Finally, `Inner` contains two more volatile fields, called `count` and `status`. These two fields are used to coordinate the rebuilding process. The `count` field contains the number of updates that occurred in the subtree rooted at this inner node since this inner node was created. The `status` field contains an integer and two booleans. When the inner node is created, it is $(0, \bot, \bot)$. It changes to a non-zero value when this inner node must be replaced with its copy during a rebuilding operation (more on that later).

The `Rebuild` data type is used when a subtree needs to be rebuilt. It has a field called `target`, which contains a pointer to the root of the subtree that needs to be rebuilt. It also has the field `parent`, which contains a pointer to the `Rebuild` node's parent, and the `index` of the `target` in the `parent` node's array of pointers to children. These two fields are sufficient to replace the `Rebuild` in the parent once the rebuild operation completes. The `status` field and the `Rebuild` type are explained in more detail later.

To work correctly, the operations of the C-IST must maintain the following invariants.

**Invariant 1** (Key presence). *For any key* k *reachable in the IST* I, *there exists exactly one path of the form* $I \xrightarrow{root} n_0 \xrightarrow{children[i_0]} (n_1 \mid r_1 \xrightarrow{target} n_1) \xrightarrow{children[i_1]} \dots \xrightarrow{children[i_{m-1}]} (n_m \mid r_m \xrightarrow{target} n_m)$, *where* $n_m$ *holds the key* k, $r_m$ *is a* Rebuild *node, and* | *is a choice between two patterns.*

**Definition 2.1** (Cover). *A root node* n *covers the interval* $\langle -\infty, \infty \rangle$. *Given an inner node* n *of degree* d *that covers the interval* $[a, b\rangle$, *and holds the keys* $k_0, k_1, \dots, k_{d-2}$ *in its* keys *array, its child* n.children[i] *covers the interval* $[k_{i-1}, k_i\rangle$, *where we define* $k_{-1} = a$ *and* $k_{d-1} = b$.
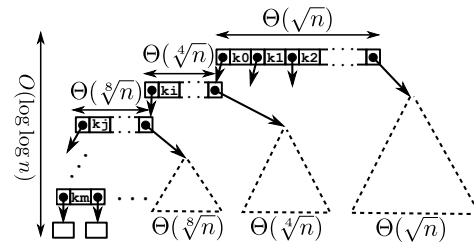
**Definition 2.2** (Content). *A node* n *contains a key* k *if and only if the path from the root of the IST* I *to the leaf with the key* k *contains the node* n. *An IST* I *contains a key* k *if and only if the* root *contains the key* k.

**Invariant 2** (Search tree). *If a node* n *covers* $[a, b\rangle$ *and contains a key* k, *then* $k \in [a, b\rangle$.

**Invariant 3** (Acyclicity). *There are no cycles in the interpolation search tree.*

**Definition 2.3** (Has-key). *Relation* hasKey$(I, k)$ *holds if and only if* I *satisfies the invariants, and contains the key* k.

In the C-IST data structure (and in the traditional, sequential IST as well), the degree $d$ of a node with cardinality $n$ is $\Theta(\sqrt{n})$. An example of an C-IST is shown in the figure below. The root has a degree $\Theta(\sqrt{n})$, its children have the degree $\Theta(\sqrt[4]{n})$, its grandchildren have the degree $\Theta(\sqrt[8]{n})$ and so on.



An *ideal C-IST* is a C-IST that is perfectly balanced. In such a C-IST, the degree of a node with cardinality $n$ is either $\lfloor \sqrt{n} \rfloor$ or $\lceil \sqrt{n} \rceil$, and the number of keys in each of the node's child subtrees is either $\lfloor \sqrt{n} \rfloor$ or $\lceil \sqrt{n} \rceil$. This property ensures that depth is bound by $O(\log \log n)$.

The interpolation search tree will generally not be ideal after a sequence of insertion and deletion operations, but its subtrees are ideal ISTs immediately after they get rebuilt.

### 2.2 Insertion and Deletion

An insertion operation first searches the tree to find a `Single` or `Empty` node that is in a location corresponding to the input key. Then, the resulting node is replaced with either one or two new nodes. An `Empty` node gets replaced with a new `Single` node that contains the new key. A `Single` node gets replaced either with a new `Single` node in case of finding an exact key match, or otherwise with a new inner node that contains two single child nodes, which contain the existing key and the newly added key.

To be able to decide when to rebalance a subtree, C-IST must track the amount of imbalance in each subtree. To achieve this, whenever a key is inserted or deleted at some leaf, the sequential IST increments the value of the `count` field, for all the inner nodes on the path from the root to the leaf that was affected [34]. Once the `count` field reaches a threshold at some inner node, the subtree below that inner node gets rebuilt. The concurrent variant, C-IST, uses a similar technique to track the amount of imbalance, but it avoids the contention at the root, where the contention is

```
struct IST is              struct Inner: Node is        struct Empty: Node is
  root: Node                 initSize: int
                             degree: int                struct Rebuild: Node is
struct Single: Node is       keys: KeyType[]             target: Inner
  key: KeyType               children: Node[]            parent: Inner
  val: ValType               status: [int, bool, bool]   index: int
                             count: int
```

**Figure 1.** Data Types

most relevant, with a scalable quiescently-consistent multi-counter [2].

After rebalancing gets triggered in some subtree, the subsequent modification operations in the corresponding subtree must fail, and help complete that rebalancing operation before retrying themselves. To ensure that this happens, the rebalancing operation sets the status field of all the inner nodes in the subtree that is being rebuilt. When replacing a child of the inner node, a modification operation, such as insertion, atomically checks the status field of that inner node. To achieve this, C-IST relies on an atomic double-compare-single-swap (DCSS) primitive, which takes two addresses, two corresponding expected values, and one new value as arguments, and behaves like a CAS, but succeeds only if the second address also matches the expected value for that second address. DCSS is bundled with a wait-free DCSS_READ primitive, which reads the fields that can be concurrently modified by a DCSS operation. This DCSS_READ operation is used by the C-IST operation whenever the value of an inner node's children array needs to be read.

Figure 2 shows the pseudocode of the insert operation. The delete operation is similar to insert. The main difference is that it either replaces a Single node with a new Empty node, or does not change the data structure if the key is not present. The deletion does not shrink the Inner nodes. This allows for some Empty nodes to be accumulated in the C-IST, but the rebuilding operations eventually occur, and remove such empty nodes.

### 2.3 Partial Rebuilding

When an insertion or deletion operation observes that a subtree rooted at some inner node *target* (henceforth, we call that subtree the "target subtree") has become sufficiently imbalanced, it triggers the rebuild operation on that subtree, which is shown in Figure 3. The rebuilding operation occurs in four steps. First, a thread announces its intention to rebuild the target subtree by creating a descriptor object of the type Rebuild, and it inserts this descriptor between the *target* and its *parent* with a DCSS operation. After the first step, other threads can cooperate on the rebuild operation. Second, this thread must prevent further modifications in the target subtree, so it does a preorder traversal to set a bit in the status field of each node of that target subtree. After the second step, no key can be added or removed from the target subtree. Third, the thread creates an ideal IST (rooted at a new inner node called *ideal*) using the old subtree's keys

```
 1 procedure insert(ist, key, val)
 2   path = []                          // Stack that saves the path.
 3   n = ist.root
 4   while true
 5     index = interpolationSearch(key, n)
 6     child = DCSS_READ(n.children[index])
 7     if child is Inner
 8       n = child
 9       path.push( [child, index] )
10     else if child is Empty | Single
11       r = createFrom(child, key, val)
12       result =
13         DCSS(n.children[index], child, r, n.status, [0,⊥,⊥])
14       if result == FAILED_MAIN_ADDRESS
15         continue                      // Retry from the same n.
16       else if result == FAILED_AUX_ADDRESS
17         return insert(ist, key, val)  // Retry from the root.
18       else
19         for each [n, index] in path
20           FETCH_AND_ADD(n.count, 1)
21         parent = ist.root
22         for each [n, index] in path
23           count = READ(n.count)
24           if count >= REBUILD_THRESHOLD * n.initSize
25             rebuild(n, parent, index)
26             break                      // exit for
27           parent = n
28         return true
29     else if child is Rebuild
30       helpRebuild(child)
31       return insert(ist, key, val)    // Retry from the root.
```

**Figure 2.** Insert Operation

(rooted at the target subtree, i.e. *target*). Finally, the thread uses a DCSS operation to replace the target subtree with the new ideal subtree in the parent. All other threads can compete to complete second, third and the fourth step.

### 2.4 Collaborative Rebuilding

When a lot of threads concurrently attempt to modify the C-IST, the rebuilding procedure described in Section 2.3 suffers from a scalability bottleneck. For example, since many threads compete to mark the old subtree in the markAndCount procedure, multiple threads will typically try to set the status bits of the same nodes. Similarly, multiple threads attempt to create the ideal version of the same target subtree in the createIdeal procedure from Figure 3, so part of the overall work gets duplicated. To reduce contention, a collaborative rebuilding algorithm is used, which allows threads to mark and rebuild different parts of the subtree in parallel.

```
32 procedure rebuild(node, p, i)
33   op = new Rebuild(node, p, i)
34   result = DCSS(p.children[i], node, op, p.status, [0,⊥,⊥])
35   if result == SUCCESS then helpRebuild(op)
36
37 procedure helpRebuild(op)
38   keyCount = markAndCount(op.target)
39   ideal = createIdeal(op.target, keyCount)
40   p = op.parent
41   DCSS(p.children[op.index], op, ideal, p.status, [0,⊥,⊥])
42
43 procedure markAndCount(node)
44   if node is Empty then return 0
45   if node is Single then return 1
46   if node is Rebuild then return markAndCount(op.target)
47   // node is Inner
48   CAS(node.status, [0,⊥,⊥], [0,⊥,⊤])
49   keyCount = 0
50   for index in 0 until length(node.children)
51     child = READ(node.children[index])
52     if child is Inner then
53       [count, finished, started] = READ(child.status)
54       if finished then keyCount += count
55       else keyCount += markAndCount(child)
56     else if child is Single then keyCount += 1
57   CAS(node.status, [0,⊥,⊤], [keyCount,⊤,⊤])
58   return keyCount
```

**Figure 3.** Rebuild Operation

```
struct Rebuild: Node is        struct Inner: Node is
  target: Inner                  initSize: int
  newTarget: Inner               degree: int
  parent: Inner                  keys: KeyType[]
  index: int                     children: Node[]
                                 status: [int, bool, bool]
                                 count: int
                                 nextMark: int
```

**Figure 4.** Modified Data Types for Collaborative Rebuilding

The previous rebuilding algorithm from Section 2.3 is changed in several ways to allow threads to perform rebuilding *collaboratively*. First of all, the markAndCount procedure, responsible for setting the status bits in a subtree, is replaced with a new markAndCountCollaborative procedure. In this new procedure, helper threads attempt to process different parts of the target subtree in parallel. They carefully avoid duplicating the work by picking a different initial index for traversing the children of the current node: each thread atomically increments a special field called nextMark to decide on its index. This is shown in Figure 5, in which we note the changes to the old markAndCount procedure.

Second, the call to createIdeal, inside the procedure helpRebuild in Figure 3, is replaced with a call to a new procedure createIdealCollaborative, in which a new root of the subtree is first created (whose children array initially contains only null-pointers) and then announced to the other threads, as before. For this purpose, the newTarget field is added to the Rebuild data type, as shown in Figure 4,

```
43 procedure markAndCountCollaborative(node)
   // ... same as markAndCount until line 48, but with
   // recursive calls to markAndCountCollaborative ...
49 if node.degree > COLLABORATION_THRESHOLD
50   while true
51     index = FETCH_AND_ADD(node.nextMark, 1)
52     if index >= node.degree then break
53     markAndCountCollaborative(node.children[index])
   // ... same as markAndCount from line 49, but with
   // recursive calls to markAndCountCollaborative ...
```

**Figure 5.** The markAndCountCollaborative Procedure

to store the root node of the new subtree. Each null-pointer in the new root of the subtree represents a "subtask" that a thread can perform by building the corresponding subtree (and then changing the corresponding null-pointer to point to this new subtree). Note that many of these subtasks can be performed in parallel. Until the new ideal IST is complete, the newTarget node effectively acts as a lock-free work pool of tasks that the threads compete for. Note that the aforementioned nextMark field in the Inner nodes is set to zero, so that it can be used for collaborative marking during the next rebuild operation.

The subtlety of these changes is to distribute work between the threads while at the same time retaining the lock-freedom property, which mandates that all work is done after a finite number of steps, even if some threads block. The collaborative rebuilding algorithm is shown in Figure 6.

### 2.5 Lookup

The lookup subroutine is shown in Figure 7. Similar to insert, it starts with an interpolation search, and repeats it at each inner node until reaching an Empty or a Single node. If it reaches a Single node that contains the specified key, it returns true. Otherwise, if lookup encounters an Empty node or a Single node that does *not* contain the specified key, it returns false.

If lookup encounters a Rebuild object, it can safely ignore the rebuild operation – lookup simply follows the target pointer to move to the next node, and continues traversal. Unlike the insert operation, lookup does not need to help concurrent subtree-rebuilding operations. Lookups do not need to help to ensure progress, and so they avoid the unnecessary slowdown. Apart from the use of the DCSS_READ operation, and checking for the Rebuild objects, *lookup* effectively behaves in the same way as the sequential interpolation tree search.

### 2.6 Summary

We presented the C-IST data types and operations precisely, but without lengthy explanations of the pseudocode. For a more descriptive discussion and concrete example scenarios, we refer the readers to related work [15].

```
59 procedure createIdealCollaborative(op, keyCount)
60   if keyCount < COLLABORATION_THRESHOLD then
61     newTarget = createIdeal(op.target, keyCount)
62   else
63     newTarget = new Inner(
64       initSize = keyCount,
65       degree = 0, // Will be set to final value in line 76.
66       keys = new KeyType[⌊√keyCount⌋ - 1],
67       children = new Node[⌊√keyCount⌋],
68       status = [0, ⊥, ⊥], count = 0, nextMark = 0)
69   if not CAS(op.newTarget, null, newTarget) then
70     // Subtree root was inserted by another thread.
71     newTarget = READ(op.newTarget)
72   if keyCount < COLLABORATION_THRESHOLD then
73     while true
74       index = READ(newTarget.degree)
75       if index == length(newTarget.children) then break
76       if CAS(newTarget.degree, index, index + 1) then
77         if not rebuildAndSetChild(op, keyCount, index)
78           return newTarget
79     for index in 0 until length(newTarget.children)
80       child = READ(newTarget.children[index])
81       if child == null then
82         if not rebuildAndSetChild(op, keyCount, index)
83           return newTarget
84   return newTarget
85
86 procedure rebuildAndSetChild(op, keyCount, index)
87   // Calculate the key interval for this child, and rebuild.
88   totalChildren = ⌊√keyCount⌋
89   childSize = ⌊keyCount / totalChildren⌋
90   remainder = keyCount % totalChildren
91   fromKey = childSize * index + min(index, remainder)
92   childKeyCount = childSize + (index < remainder ? 1 : 0)
93   child = createIdeal(op.target, fromKey, childKeyCount)
94   if index < length(op.newTarget.keys)
95     key = findKeyAtIndex(op.target, fromKey)
96     WRITE(op.newTarget.keys[index], key)
97   // Set new child, check if failed due to status change.
98   result = DCSS(op.newTarget.children[index],
99     null, child, op.newTarget.status, [0, ⊥, ⊥])
100  return result != FAILED_AUX_ADDRESS
```

**Figure 6.** The `createIdealCollaborative` Procedure

```
101 procedure lookup(ist, key)
102   n = ist.root
103   while true
104     if n is Inner then
105       index = interpolationSearch(key, n)
106       n = DCSS_READ(n.children[index])
107     else if n is Single then return n.k == key ? n.v : null
108     else if n is Empty then return null
109     else if n is Rebuild then n = n.target
```

**Figure 7.** Lookup Operation

## 3 Correctness

In this section, we formalize the concurrent IST, and prove that the operations are safe, linearizable and non-blocking. We begin by introducing the concept of an abstract set, and then define the correspondence between the abstract set and the concurrent IST.

**Definition 3.1** (Abstract set). An *abstract set* $\mathbb{A}$ is a mapping $\mathbb{A} : \mathbb{K} \rightarrow \{\top, \bot\}$, where $\mathbb{K}$ is the set of all keys, and which is true ($\top$) for keys that are present in the set. The abstract set operations are: lookup$(\mathbb{A}, k) = \top \Leftrightarrow k \in \mathbb{A}$, insert$(\mathbb{A}, k) = \{k' : k' \in \mathbb{A} \vee k' = k\}$, and delete$(\mathbb{A}, k) = \{k' : k' \in \mathbb{A} \wedge k' \neq k\}$.

**Definition 3.2** (Validity). An IST $\mathtt{I}$ is *valid* if and only if the Invariants 1, 2 and 3 hold.

**Definition 3.3** (Consistency). An IST $\mathtt{I}$ is *consistent* with the abstract set $\mathbb{A}$ if and only if $\forall k \in \mathbb{A} \Leftrightarrow \mathtt{hasKey}(\mathtt{I}, k)$, (as per Definition 2.3). The IST `lookup` operation is consistent with an abstract set lookup if and only if $\forall k, \mathtt{lookup}(\mathtt{I}, k) = \mathtt{lookup}(\mathbb{A}, k)$. The IST `insert` and `delete` operations are consistent with the abstract set insert and delete if and only if for all keys $k$ and initial states $\mathtt{I}$ consistent with some abstract set $\mathbb{A}$, they leave the IST in a state $\mathtt{I}'$, consistent with some abstract set $\mathbb{A}'$, such that insert$(\mathbb{A}, k) = \mathbb{A}'$ or delete$(\mathbb{A}, k) = \mathbb{A}'$, respectively.

**Theorem 3.4** (Safety). *An IST $\mathtt{I}$ that uses non-collaborative rebuilding from Figure 3 is always valid and consistent with some abstract set $\mathbb{A}$. IST operations are consistent with the abstract set semantics.*

We prove safety inductively – it is easy to see that the newly created IST is consistent with the empty abstract set. For the inductive case, we must show that after each mutation, the IST remains valid, and that it either remains consistent with the same abstract set $\mathbb{A}$, or becomes consistent with some other set $\mathbb{A}'$, according to the abstract set semantics. After that, we extend the proof to the IST variant that uses collaborative rebuilding.

**Lemma 3.5** (Freezing). *After `markAndCount` in line 38 of Figure 3 completes, all the nodes in the corresponding subtree have `status` different than $[0, \bot, \bot]$.*

*Proof.* Consider the CAS in line 48 of Figure 3. Exactly one thread will successfully execute this CAS. After this happens, none of the entries in the `children` array will be modified, since all the DCSS invocations are conditional on `status` being zero. Assuming that the statement holds inductively for the children, when a thread reaches the CAS in line 57, all the inner nodes in the corresponding subtree have a non-zero `status`.                                                   □

**Lemma 3.6** (End of Life). *If the field `status` of an inner node $\mathtt{n}$ is $[0, \bot, \bot]$ at $t_0$, then the node $\mathtt{n}$ is reachable from the root at some time $t_0$.*

*Proof.* Assume the converse – node $\mathtt{n}$ is unreachable at $t_0$ and it has `status` set to $[0, \bot, \bot]$. The only instruction that makes inner nodes unreachable is the DCSS in line 41 of Figure 3, which is executed on $\mathtt{n}$'s parent. This instruction is executed after `markAndCount` in line 38 completes, so the inner nodes in the corresponding subtree must have a `status` value different than $[0, \bot, \bot]$, by Lemma 3.5.                □

**Lemma 3.7** (Cover stability). *If a node n covers the interval $[a, b\rangle$ at some time $t_0$ (as per Definition 2.1), then n covers the same interval $\forall t > t_0$, as long as n is reachable at the time t.*

*Proof.* We prove this inductively. For the base case, the statement trivially holds for the root. Next, assume that it holds for the parent of some node n. If n is reachable at some time $t$, then so is its parent, whose cover did not change. Since the cover of n is defined based on the cover of its parent and the contents of the immutable keys array (as specified in Definition 2.1), the cover of n also does not change after its creation. □

*Proof of Theorem 3.4.* We consider all execution steps that mutate the data structure, and reason about the respective change in the corresponding abstract set.

Consider the DCSS in line 13 of Figure 2. If this DCSS succeeds at $t_0$, then the status field of the corresponding node must be $[0, \perp, \perp]$, and the node is reachable at $t_0$, by Lemma 3.6. Before $t_0$, the tree is inductively valid and consistent with some set $\mathbb{A}$. After $t_0$, the new node (created in the createFrom subroutine) becomes reachable from the root. The new node contains the new key, and (if child was non-Empty, and its key different than the new key) the old key from child. Therefore, this change is consistent with the abstract set semantics. Furthermore, Invariants 1, and 3 are trivially preserved. To see that Invariant 2 holds, note that all the nodes on the path from the root to node are reachable, since node is reachable by Lemma 3.6. Each of those nodes covers exactly the same interval that it covered when interpolationSearch ran in line 5 of Figure 2, by Lemma 3.7. The node that was created by the createFrom subroutine also respects Invariant 2. Therefore, the new key is inside the intervals covered by those nodes, and the search tree property is preserved.

Next, consider the DCSS in line 41 of Figure 3. After the markAndCount call in line 38 completes, that DCSS cannot succeed, by Lemma 3.5. Thus, there will be no further updates of the children arrays in that subtree, so the corresponding subtree is effectively immutable. All createIdeal calls occur after the subtree becomes immutable, and each such call will create a new subtree that corresponds to the same abstract set as the old subtree, so the DCSS in line 41 does not change the consistency or the validity of the IST.

It is easy to see that the DCSS in line 34 of Figure 3 changes neither the validity nor the consistency of an IST, since it only inserts up to one Rebuild node between some two Inner nodes.

Finally, neither the CAS instructions in lines 48 and 57 of Figure 3, nor the FETCH_AND_ADD in line 20 of Figure 2, affect validity or consistency of the IST, since these properties are independent of the status and the count fields. □

Next, we need to show safety for the collaborative IST variant. The crux of the proof is to show that the procedure markAndCountCollaborative and the procedure createIdealCollaborative have the same semantics as their non-collaborative variants.

**Lemma 3.8** (Collaborative Freezing). *After the procedure markAndCountCollaborative from Figure 5 completes, all the nodes in the corresponding subtree have status different than $[0, \perp, \perp]$.*

*Proof.* After lines 49-53 of markAndCountCollaborative complete, some subset of the nodes below the target node will have their status different than $[0, \perp, \perp]$, due the CASes in lines 48 and 57. This is not necessarily for all of the nodes below the target node, since some threads might still be processing their subtrees at the time when the current thread exits the loop in lines 49-53. Nevertheless, the loop from Figure 3 in lines 50-55, which also appears in markAndCountCollaborative (after the loop in lines 49-53), ensures that each descendant node has their status set to a non-$[0, \perp, \perp]$ value before the procedure ends, so the proof is similar to that of Lemma 3.5. □

**Corollary 3.9** (Collaborative End of Life). *Lemma 3.6 holds when the markAndCount procedure is replaced with the call to markAndCountCollaborative.*

**Lemma 3.10** (Collaborative Ideal Tree Creation). *Let T be a subtree below the target node, which corresponds to some abstract set $\mathbb{S}$, and on which markAndCountCollaborative completed. Once the createIdealCollaborative completes, either the tree T gets replaced with a new ideal tree T' that corresponds to the same abstract set $\mathbb{S}$, or the status of the target node changed to a non-$[0, \perp, \perp]$ value.*

*Proof.* First, not that the tree T is effectively-immutable by Lemma 3.8 and Corollary 3.9, so the entries in its children arrays never change.

The case in which the keyCount is below the threshold creates an already-complete ideal tree newTarget. In the other case, local variable newTarget is set to an empty root node for the new ideal tree. If the CAS in line 69 succeeds, then op.newTarget points to this new root. Conversely, if that CAS fails, then op.newTarget points to a new root created by another thread, which, due to the effectively-immutable T based on which it was created, must be structurally identical to the node in the local variable newTarget (but is a different memory object in the heap). Thus, all threads that move beyond the line 71, have the pointer to the same root node in their newTarget variable, which corresponds to the pointer in the op.newTarget field.

Next, consider the case in which the status field of the op.target node never changed to a non-$[0, \perp, \perp]$ value. In this case, the DCSS in line 98 of the rebuildAndSetChild procedure can never fail due to a change in the auxiliary

address. Therefore, in this case, `rebuildAndSetChild` always returns `true`, so the `createIdealCollaborative` cannot return early in lines 78 or 83, and both loops must finish before `createIdealCollaborative` returns.

Consider the loop in lines 73-78. After any thread exits this loop, an entry $i$ of the `newTarget.children` array contains either `null`, or the $i$-th child of the ideal tree corresponding to $\mathbb{S}$. Moreover, after any thread exits this loop, it must be true that `newTarget.degree` is equal to the length of the `newTarget.children` array.

Next, consider the loop in lines 79-83. After any thread exits this loop, it must be true that every entry $i$ of the `newTarget.children` array contains the $i$-th child of the ideal tree corresponding to $\mathbb{S}$. Moreover, $i$-th entry of `newTarget.keys` array contains the $i$-th key of the ideal tree corresponding to $\mathbb{S}$ (due to the write in line 96).

Therefore, if the `status` field of the `op.target` node did not change to a non-[0, ⊥, ⊥] value, then the procedure `createIdealCollaborative` returns an ideal tree corresponding to $\mathbb{S}$.

Finally, consider the case in which the `status` field of the `op.target` node did change to a non-[0, ⊥, ⊥] value. Here, the claim trivially holds. □

**Theorem 3.11** (Collaborative Safety). *An IST I that uses collaborative rebuilding from Figures 5 and 6 is always valid and consistent with some abstract set $\mathbb{A}$. IST operations are consistent with the abstract set semantics.*

*Proof.* Similar to the proof for Theorem 3.11, using Lemma 3.8 and Lemma 3.10 to show that the `DCSS` in line 41 does not cause a change in the abstract set.

In particular, if the `createIdealCollaborative` returns prematurely due to a non-[0, ⊥, ⊥] value in the `status` field, then the `DCSS` in line 34 in Figure 3 will fail and cause a restart from the root, due to a rebuild operation in one of the ancestors. □

Linearizability follows naturally from the safety proofs, since all the state-changing execution steps occur at the atomic operations.

**Corollary 3.12** (Linearizability). *Lookup, insertion and deletion operations are linearizable, in both the non-collaborative and collaborative variant of the concurrent IST data structure.*

*Proof.* For insertion and deletion, this follows immediately from the proofs of Theorems 3.4 and 3.11, since the `DCSS` in line 13 is the linearization point. The linearization point of a lookup is the last `DCSS_READ` in line 106 of Figure 7 that was done on a node whose `status` field was [0, ⊥, ⊥]. □

We first state the non-blocking properties of the non-collaborative IST variant. After that, we extend the proof to the collaborative variant.

**Theorem 3.13** (Non-Blocking). *In the non-collaborative IST variant, insertion and deletion are lock-free, and lookup is wait-free.*

*Proof.* By Invariant 3, every path in the IST is finite. The `lookup` subroutine in Figure 7 executes a finite number of steps in each loop iteration before advancing to the next node on the path. At any point, there can be up to $p$ insertion threads that are racing to extend this path. The `lookup` operation can potentially race with these inserting threads (for example, if they are all repetitively extending the same leaf node). However, after a finite number of steps, each of these $p$ threads will increment the `count` field beyond a threshold and trigger the rebuild of the enclosing subtree. At that point, the path can no longer be extended, and the `lookup` observes a `Rebuild` node, and terminates in a finite number of steps. Therefore, `lookup` is wait-free.

For the insertion and deletion, identify all the instructions that do not change the abstract state $\mathbb{A}$, but can fail an instruction that is a linearization point. We call these instructions *housekeeping*. The `DCSS` in line 34 of Figure 3 and the `CAS` in line 48 of Figure 3, are the only such housekeeping instructions. These instructions can be executed only finitely many times before some instruction changes $\mathbb{A}$.

To see this, define the *pending count* of an inner node `n` as the value of its field `count` increased by the number of threads $p$ that are currently preparing to increment `n.count` (note: each such thread can increment `n.count` at most once without changing $\mathbb{A}$). Next, identify all nodes whose pending count crossed the threshold for rebuilding, and call them *triggered*. Define an abstract variable $\mathbb{V}_0$ as the triggered nodes that are directly pointed to by the `children` array. Define $\mathbb{V}_1$ as the total number of inner nodes in the subtrees of all triggered nodes whose `status` field is 0. Note that $\mathbb{V} = \mathbb{V}_0 + \mathbb{V}_1$ is finite. Furthermore, as long as $\mathbb{A}$ stays the same, housekeeping instructions can only strictly decrease $\mathbb{V}$. When $\mathbb{V}$ is 0, neither the `DCSS` in line 34 of Figure 3 nor the `CAS` in line 48 of Figure 3, can succeed, so they cannot fail an instruction that would cause a linearization point. To complete the proof, note that two invocations of the `DCSS` in line 13 of Figure 2 are distanced by a finite number of steps, by the same argument as for the `READ`s in the `lookup` subroutine. □

In the collaborative variant, the `lookup` operation is wait-free by the same arguments. To prove that insertion and deletion are also non-blocking, we need to show that the `markAndCountCollaborative` procedure, as well as the `createIdealCollaborative` procedure, both complete in a finite number of steps.

**Lemma 3.14.** *The `markAndCountCollaborative` procedure and the `createIdealCollaborative` procedure are wait-free.*

*Proof.* First, it is easy to see that helper procedures, such as `createIdeal`, finish in a finite number of steps (because the tree is finite).

Next, note that the loops in lines 50-55 of Figure 3, and the lines 79-83 of Figure 6, both make progress in each iteration, and finish within a finite number of execution steps.

The same is true of the loop in lines 50-53 of Figure 5 – after every iteration of the loop, the `nextMark` field is strictly increased. Thus, `markAndCountCollaborative` is wait-free.

The loop in lines 73-53 of Figure 6 contains a `CAS` on the `degree` field that makes progress when successful. However, that `CAS` fails only if another thread successfully executes the same `CAS` on the `degree` field, incrementing it by 1. Therefore, an iteration of that loop makes progress in both cases, and completes in a finite number of iterations. This proves that the `createIdealCollaborative` procedure is wait-free. □

**Corollary 3.15.** *In the collaborative IST variant, insertion and deletion are lock-free, and lookup is wait-free.*

*Proof.* Follows from Theorem 3.13 and Lemma 3.14. □

## 4 Complexity Analysis

The complexity proof for C-IST follows the argument for sequential ISTs [34], with a few modifications – most notably, it accounts for the fact that at any time there can be up to $p$ threads that are concurrently modifying the C-IST.

**Lemma 4.1.** *Let $p$ be the number of concurrent threads that are modifying a C-IST. Worst-case depth of a C-IST that contains $n$ keys is $O(p + \log n)$.*

*Proof.* Let u be some node and v be its parent. In the worst case, all keys that were updated below the node v, were also below its child u. In our case, the rebuild threshold $R = 1/4$, so there were at most $v.initSize/4 + p_v$ such updates, where $v.initSize$ is the initial size of the node v, and $p_v$ is the number of concurrent threads that are currently modifying the subtree below v. We add $p_v$ because those threads possibly inserted the new keys, but did not yet update the counts on the path to the root.

Since u was an ideal tree at the time when its parent v was created, we have that:

$$|u| \leq \sqrt{v.initSize} + v.initSize/4 + p_v \leq v.initSize/2 + p_v \quad (1)$$

Next, consider the longest path $L$ from the root r to some leaf. On this path, the total number of concurrent updates that have not yet updated the `count` fields is $p_L \leq \sum_v p_v = p$. The new nodes inserted by these $p_L \leq p$ updates, which we denote by $w_1, w_2, \ldots, w_{p_L}$ must be at the suffix of the path, that is:

$$L = r \to v_1 \to \ldots \to v_{d_0} \to w_1 \to \ldots \to w_{p_L} \quad (2)$$

Otherwise, if the nodes represented by those updates were moved elsewhere, then a rebuild must have been triggered, meaning that the `count` fields *had already been updated*, which is a contradiction with how we defined $w_1, \ldots, w_{p_L}$.

Therefore, these $p_L \leq p$ updates can only increase the depth by $p_L$. Let us ignore them in the size consideration momentarily. For the prefix $r \to v_1 \to \ldots \to v_{d_0}$ of the path $L$, we have:

$$r.initSize \geq 2 \cdot |v_1| \geq \ldots \geq 2^{d_0} \cdot 1 \quad (3)$$

Furthermore, $r.initSize \leq |r| \cdot (1 - R)^{-1}$ so the depth $d$ of the path $L$ is:

$$d = p_L + d_0 \leq p + \log_2(r.initSize) \leq O(p + \log |r|) \quad (4)$$

which proves the claim of the lemma. □

**Lemma 4.2.** *The worst-case amortized cost of insert and delete operations, without including the cost of searching for the node in the C-IST, is $O(\gamma(p + \log n))$, where $\gamma$ is a bound on average interval contention.*

*Proof.* Let us consider the execution at a node $v$ between two rebuilding phases. We can split the work performed at the node itself and by operations which traverse it into three disjoint categories:

- Work performed to rebuild the node.
- Work by operations which traverse the node $v$ but *fail* in line 17 and restart from the root.
- Work performed by *successful* operations, which do not restart from the root.

First, note that the amortized cost of a single successful operation is no worse than the worst-case length of a root-to-leaf path, which is $O(p + \log n)$, by Lemma 4.1. Second, notice that, by standard non-blocking amortization argument, a successful operation can be on average responsible for at most $\gamma$ distinct *failed* operations, each of which has worst-case cost $O(p + \log n)$. Third, notice that the work necessary to rebuild the subtree rooted at node $v$ is in the order of $nodes(v) \leq (1 + R) \cdot v.initSize + p$. However, this rebuild process is triggered only when at least $R \cdot v.initSize$ operations have *succeeded* in the subtree, as only successful operations increment the counter and therefore trigger a rebuild.

We therefore obtain that the total amortized cost of a successful operation is upper bounded by

$$O(\gamma(p + \log n)) + ((1 + R) \cdot v.initSize + p)/R \cdot v.initSize)$$
$$= O(\gamma(p + \log n)).$$

□

To analyze the expected amortized cost of operations, we introduce a concept of a $\mu$-random C-IST, identically to the sequential variant.

**Definition 4.3.** Let $\mu$ be a probability density function on reals. A $\mu$-random C-IST of size $n$ is generated as follows:

1. Independently draw $n_s$ real keys according to the probability density $\mu$, and use them to create an ideal IST.
2. Do a sequence of updates $U_1, \ldots, U_m$, with $i$ insertions and $d$ deletions, so that $m = i + d$ and $n = n_s + i - d$. An insertion picks a random key according to the probability density $\mu$. A deletion picks a random key that is present in the C-IST, according to a uniform probability density.

We will now establish some bounds on the expected size of subtrees in C-IST. We will track two components of the C-IST size – the size contribution due to the updates that are already completed, and the size contribution due to slow updates due to straggler threads that have inserted a node, but have not yet incremented the count fields, nor checked whether a rebuild is required. We will assume a standard non-malicious *oblivious* shared-memory scheduler, whose actions are independent of the input distribution and of the threads random coin flips.

**Definition 4.4.** Let the *initial size* of a node denote the node count of its corresponding subtree at the point where it was last rebuilt. Let the *stable size* denote the number of keys that were inserted or deleted, and for which the count fields of the nodes on the respective paths were updated, up to the current node, at the current moment in time.

We show that, in a $\mu$-random C-IST, as long as nodes have size $\Omega(p)$, the stable size distributes evenly over the subtrees.

**Lemma 4.5.** *Let $\mu$ be a density function with finite support $[a, b]$, let $v$ be the root of a $\mu$-random C-IST of initial size $n_0$, such that $n_0 R > p$, of stable size $n$ at the current point. Let $u$ be a child of the node $v$. Then, the stable size of the subtree rooted at $u$ will be $O(n^{1/2})$, at any future point in time until the next rebuild of node $u$ with probability $\geq 1 - O(n^{-1/2})$.*

*Proof.* We start by noticing that, by the rebuild conditions, and since there can be at most $p - 1$ pending operations at any given point before the next rebuild, we have that $n \leq (1 + R)n_0 + p \leq (1 + 2R)n_0$. Also, we know that there could have been at most $i \leq Rn_0$ inserts into $v$ since it was last rebuilt.

At the time when $v$ was rebuilt, its child $u$ had the size $\sqrt{n_0} \pm 1$. Let $X$ be a random variable that denotes the number of elements that were stored in $u$ since that time. Using the result from Theorem A from sequential ISTs [34], the expectation and the variance of $X$ can be bounded as follows:

$$E(X) = \sqrt{n_0} \cdot \frac{i}{n_0 + 1} \leq R \cdot \sqrt{n_0}. \tag{5}$$

$$Var(X) \leq E(X) \cdot (1 + \frac{i}{n_0 + 1}) \leq 2 \cdot R \cdot \sqrt{n_0}. \tag{6}$$

We can therefore apply Chebyshev's inequality $P(|X - E(X)| \geq t) \leq \frac{Var(x)}{t^2}$ for $t = (1 + R)n_0^{1/2}$, to bound the probability that $X$ exceeds the square root size by a factor of $1 + R$:

$$P(X \geq (1 + R) \cdot n_0^{1/2}) \leq \frac{2R}{(1 + R)^2} \cdot n_0^{-1/2}. \tag{7}$$

Thus, the claim about the stable size follows.

$\square$

**Lemma 4.6.** *Let $\mu$ be a probability density with a finite support $[a, b]$. The expected total cost of processing a sequence of $n$ $\mu$-random insertions and uniformly random deletions into an initially empty C-IST is $O(n(\log \log n + p)\gamma)$, where $\gamma$ is a bound on average interval contention.*

*Proof.* We proceed by defining the following token generation scheme for *successful* operations: we say that a traversal corresponding to an insert or remove operation generates a token at a node $v$ if it passes the node on its walk from the root to a leaf. (Notice that this also counts the number of fetch-and-increment operations the walk will generate on its backward path.) We can split such generated tokens into *internal tokens*, which are generated at internal nodes, and *leaf tokens*, which are generated at leaves. Notice that the only way in which a walk can generate multiple leaf tokens is if it encounters additional concurrent operations.

We proceed to first bound the expected number of internal tokens generated by a successful root-to-leaf walk. If we denote this number of additional internal tokens generated by the successful update when starting from a node of stable size $m$ by $T(m)$, we have that, by Lemma 4.5 and Lemma 4.1,

$$T(m) \leq 1 + T(O(m^{1/2})) + O(m^{-1/2}) \cdot O(\log m), \tag{8}$$

where the recursion has to stop at nodes of maximum initial size $p$. However, the number of tokens generated nodes below this size is at most $p$, and therefore, by the above recursion, we obtain that the expected number of generated internal tokens by a traversal can be at most $O(p + \log \log n)$. We are therefore left with the task of bounding the number of *leaf tokens* generated by the successful walk. Since the operations are non-blocking, the average number of such tokens per walk is $O(p)$, which yields a total average token bound of $O(p + \log \log n)$. Similarly to Lemma 4.2, we notice that each token corresponds to a constant amount of work by some operation. Therefore, the amortized cost of successful operations is $O(p + \log \log n)$.

Our argument so far has only taken into account successful traversals, which do not restart from the root. If we simply amortize each failed traversal against the successful one which caused it to fail, we obtain that each successful traversal $i$ can be charged for $O(\gamma_i(p + \log \log n))$ tokens,

where $\gamma_i$ is its interval contention. Therefore, we get that the amortized expected cost of an operation is $O(\gamma(p+\log \log n))$, as claimed.

We note that stronger bounds on amortized expected cost could potentially be obtained by attempting to leverage the structure of the distribution to bound the probability of a collision. However, this is not immediate, since the smoothness property does not imply a small upper bound on the probability of a key. □

**Lemma 4.7.** *Let $\mu$ be a smooth probability density, as defined by Mehlhorn and Tsakalidis [34], for a parameter $\alpha$, such that $\frac{1}{2} \leq \alpha < 1$, and let $v$ be the root of a $\mu$-random C-IST with parameter $\alpha$. Then, the expected search time in the* `v.children` *is $O(1)$.*

*Proof.* Since the `children` array is immutable, the proof is precisely identical to the one in the related work on sequential ISTs [34]. □

**Lemma 4.8.** *Let $\mu$ be a smooth probability density for a parameter $\alpha$, such that $\frac{1}{2} \leq \alpha < 1$. The expected search time in a $\mu$-random IST of size $n$ is $O(\log \log n + p)$.*

*Proof.* The proof follows from Lemma 4.7 and the argument in Lemma 4.6, and the bound of $\leq p$ on the size of a root extension before a rebuild is triggered from Theorem 3.13. □

# 5 Evaluation

## 5.1 Basic Experimental Evaluation

This section summarizes the experimental setup, and the results reported in the related work that proposes and describes C-IST in detail [15].

We implemented the concurrent IST in C++, and we ran the benchmarks on a NUMA system with four Intel Xeon Platinum 8160 3.7GHz CPUs, each of which has 24 cores and 48 hardware threads. Within each CPU, cores share a 33MB LLC, and cores on different CPUs do not share any caches. The system has 384GB of RAM, and runs Ubuntu Linux 18.04.1 LTS. Our code was compiled with GCC 7.4.0-1, with the highest optimization level (`-03`). Threads were *pinned* to cores such that thread counts up to 48 ran on only one CPU, thread counts up to 96 run on only two CPUs, and so on. We used the fast scalable allocator jemalloc 5.0.1-25. When a memory page is allocated on our 4-CPU Xeon system, it has an *affinity* for a single CPU, and other CPUs pay a penalty to access it. We used the `numactl -interleave=all` option to ensure that pages are evenly distributed across CPUs.

The goal of this section is to show that the amortized $O(\log \log n)$ running time induces performance improvements on datasets that are reasonably large. We therefore evaluate the C-IST operations against other comparable data structures in Section 5.1.1, where we show, for 1 billion keys,

improvements ranging from 15-50% compared to the $(a, b)$-tree [13] (the next best alternative), depending on the ratio of updates and lookups. To further characterize the performance, we compare the average key depth and the impact on cache behavior in Section 5.1.2. We conclude with a comparison of memory footprints in Section 5.1.3.

### 5.1.1 Comparison of the Basic Operations

Figure 10 shows the throughput of concurrent IST operations, compared against other sorted set data structures, for dataset sizes of $k = 2 \cdot 10^8$ and $k = 2 \cdot 10^9$ keys, and for $u = 0\%$, $u = 1\%$, $u = 10\%$ and $u = 40\%$, where $u$ is the ratio of update operations among all operations.

In all cases, C-IST operations have much higher throughput than Natarajan and Mittal's non-blocking binary search tree (NM), and concurrent AVL trees due to Bronson (BCCO). For update ratios $u = 0\%$ and $u = 1\%$, concurrent IST also has a higher throughput compared to Brown's non-blocking $(a, b)$-tree. The underlying cause for better throughput is a lower rate of LLC misses due to IST's doubly-logarithmic depth. For higher update ratios $u = 10\%$ and $u = 40\%$, the cost of concurrent rebuilds starts to dominate the gains of doubly-logarithmic searches, and ABTree has a better throughput for $k = 2 \cdot 10^8$ keys. Above $k = 2 \cdot 10^9$ keys, ISTree outperforms ABTree even for the update ratio of $u = 40\%$.

### 5.1.2 Average Depth and Cache Behavior

The main benefit of C-IST's expected-$O(\log \log n)$ depth is that the key-search results in less cache misses compared to other tree data structures. The plot shown below compares the average number of pointer hops required to reach a key (error bars show min/max values over all trials), for dataset sizes from $2 \cdot 10^6$ to $2 \cdot 10^9$ keys. While the average depth is 20-40 for NM and BCCO, the average ABTree depth is between 6 and 10, and the average C-IST depth is below 5.



The differences in average depths between these data structures correlate with the average number of cache misses. Figure 11 compares the average number of last-level cache-misses between the different data structures, for different update ratios $u$. For the dataset size of $2 \cdot 10^9$ keys, ISTree operations undergo 2× less cache misses, and slightly fewer cache misses than ABTree.

### 5.1.3 Memory Footprint

Due to using a lower number of nodes for the same dataset, the average space overhead is lower for the C-IST than the other data structures. Figure 8 shows the different memory footprints for four different dataset sizes. C-IST has a relative space overhead of ≈30-100%, whereas the overhead of the other data structures is between ≈120-400%.
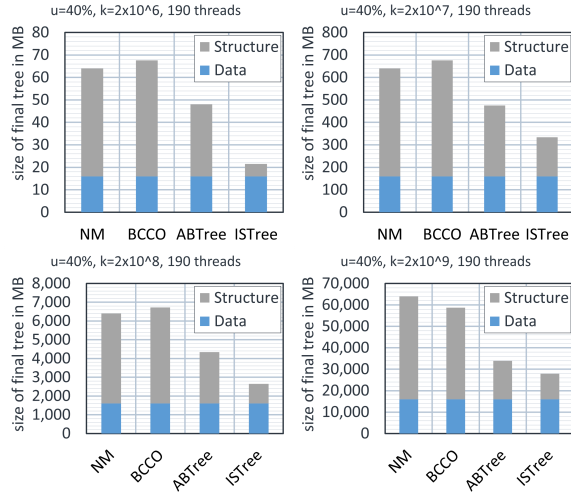
**Figure 8.** Memory Footprint Comparison

## 5.2 Extended Experimental Evaluation

In this section, we show additional experimental data that did not fit into the main body of the paper.

**Excluding data structures from graphs.** As we mentioned in Section 5.1, we also compared ISTree with many other concurrent search trees, including the external lock-free BST of Ellen et al. [22] (EFRB), the internal lock-free BST of Howley and Jones [30] (HJ), the partially external lock-free BST of Drachsler et al. [21] (DVY), the internal lock-free BST of Ramachandran and Mittal [55] (RM), the external lock-based BST of David et al. [18] (DGT), and the external lock-free BST (BER), lock-free relaxed AVL tree (RAVL), and lock-free Chromatic tree (Chromatic) of Brown et al. [13].

Unfortunately, the implementations of HJ and RM (which we obtained from ASCYLIB [18]) failed basic validation checksums. Even after attempting to fix the implementations as much as possible, the failures were persistent. Since the implementations are not correct, it does not make sense to benchmark them. The remaining implementations, EFRB, DVY, DGT, DER, RAVL and Chromatic were all consistently outperformed by BCCO and, hence, merely obscured the results. To give a sense of what the full data looks like, Figure 9 compares the implementations that do not fail checksum validation using the same experimental setup as Section 5.1.1. (RAVL and Chromatic perform similarly to BER and are excluded.) Thus, we decided not to include any of these data structures in our graphs.

**Throughput.** Figure 10 shows the throughput of the different concurrent data structures for different update percentages (0%, 1%, 10% and 40%), and dataset sizes ranging from $10^6$ to $10^9$.

**LLC misses.** Figure 11 shows the average count of last-level cache-misses for different update percentages (0%, 1%, 10% and 40%), and dataset sizes ranging from $10^6$ to $10^9$.



**Figure 9.** Full comparison of basic operations for all (non-failing) data structures, higher is better

**Collaborative rebuilding.** Figures 12 and 13 show the breakdown of the execution time, with non-collaborative and collaborative rebuilding, respectively. Plots in Figure 12 show that, as the ratio of updates grows beyond 10%, the execution time quickly becomes dominated by rebuilding. In fact, most of the rebuilding time is spent in helping operations in which a created subtree ends up being discarded. The proportion of discarded subtrees grows with the number of threads. Plots in Figure 13 show that, when using collaborative rebuilding, almost no time is spent in discarding the subtrees.

**Using a multicounter in the root node.** Our concurrent IST implementation uses atomic fetch-and-add instructions to implement update counters in inner nodes, but uses a multicounter in the root node to reduce contention. Figure 14 compares the throughput of a concurrent IST with a multicounter against an IST with a simple fetch-and-add-counter in the root node, when using 190 threads. In the variant that uses a simple fetch-and-add-counter in the root, throughput is reduced to $15\% - 33\%$ compared to the multicounter variant.

**No-SQL database workload.** We study a simple *in-memory database management system* called DBx1000 [59], which is used in multi-core database research. DBx implements a simple relational database, which contains one or more
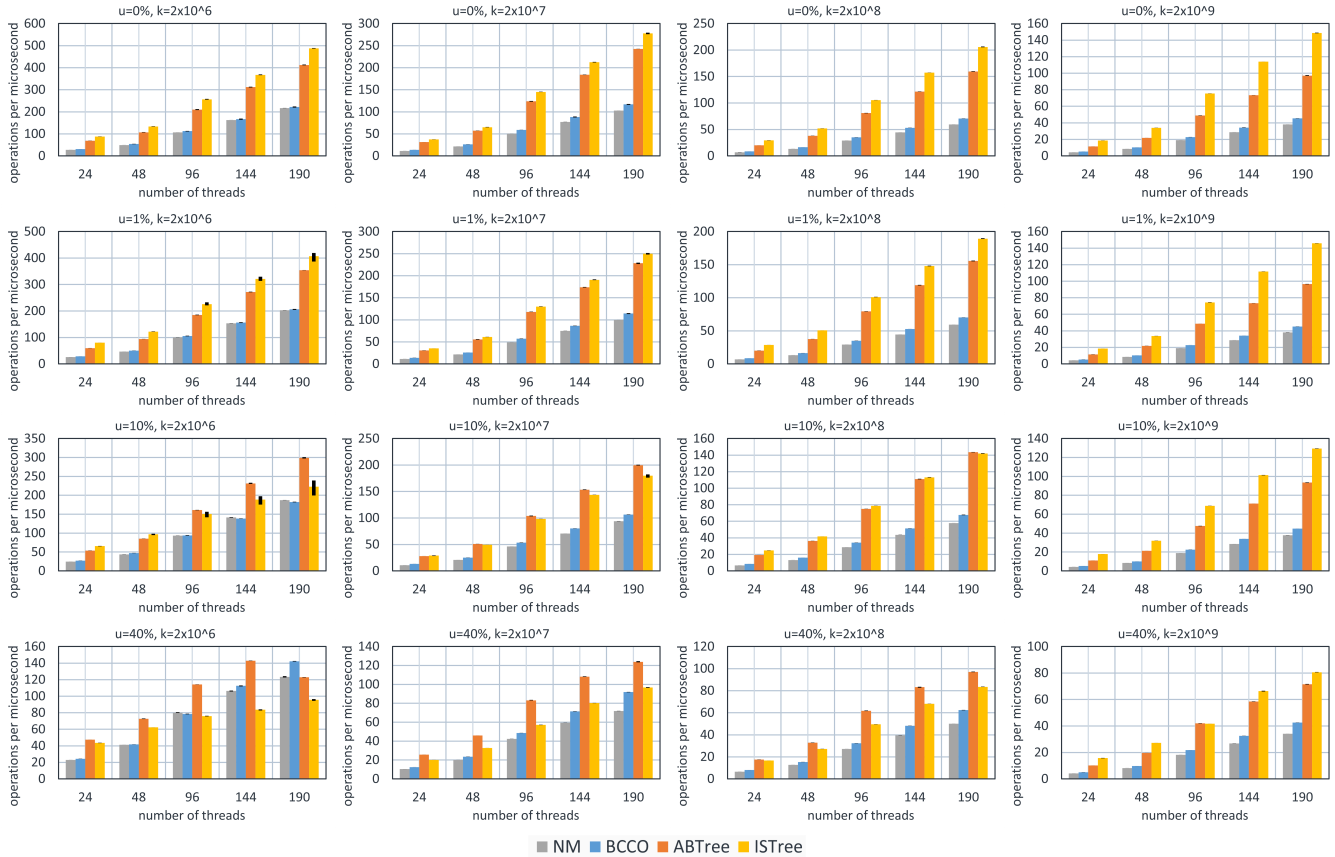
Aleksandar Prokopec, Trevor Brown, and Dan Alistarh



**Figure 10.** Comparison of Basic Operations, Higher is Better

*tables.* Each table can have one or more *key fields* and associated *indexes.* Each index allows processes to query a specific key field, quickly locating any rows in which the key field contains a desired value. We replace the default index implementation in DBx with each of the BSTs that we study.

Following the approaches in [6, 59], we run a subset of the well known Yahoo! Cloud Serving Benchmark (YCSB) core with a single table containing 100 million rows, and a single index. Each thread performs a fixed number of transactions (100,000 in our runs), and the execution terminates when the first thread finishes performing its transactions. Each transaction accesses 16 different rows in the table, which are determined by index lookups on randomly generated keys. Each row is read with probability 0.9 and written with probability 0.1. The keys that a transaction will *access* are generated according to a **Zipfian** distribution following the approach in [26].

The results, which appear in Figure 15, show how performance degrades as the distribution of *accesses* to keys becomes highly skewed. (Higher $\theta$ values imply a more extreme skew. A $\theta$ value of 0.9 is *extremely* skewed.)

**Trees containing Zipfian-distributed keys.** Since the performance of the ISTree can theoretically degrade when the

tree contains a highly skewed set of keys, we construct a synthetic benchmark to study such scenarios. In this benchmark, $n$ threads access a single instance of the ISTree, and there is a *prefilling* phase followed by a *measured* phase. In the prefilling phase, each thread repeatedly generates a key from a Zipfian distribution ($\theta = 0.5$) over the key range $[1, 10^8]$ (picking one of 100 million *possible* keys), and inserts this key into the data structure (if it is not already present). This continues until the data structure contains 10 million keys (only 10% of the key range), at which point the prefilling phase ends. In the measured phase, all threads perform $u\%$ updates and $(100 - u)\%$ searches (for $u \in \{0, 1, 10\}$) on keys drawn from the same Zipfian distribution, for 30 seconds. This entire process is repeated for multiple trials, and for thread counts $n \in \{24, 48, 96, 144, 190\}$ (with at least one core left idle to run system processes). The results, which appear in Figure 16, suggest that the ISTree can be surprisingly robust in scenarios where it contains a highly skewed distribution.
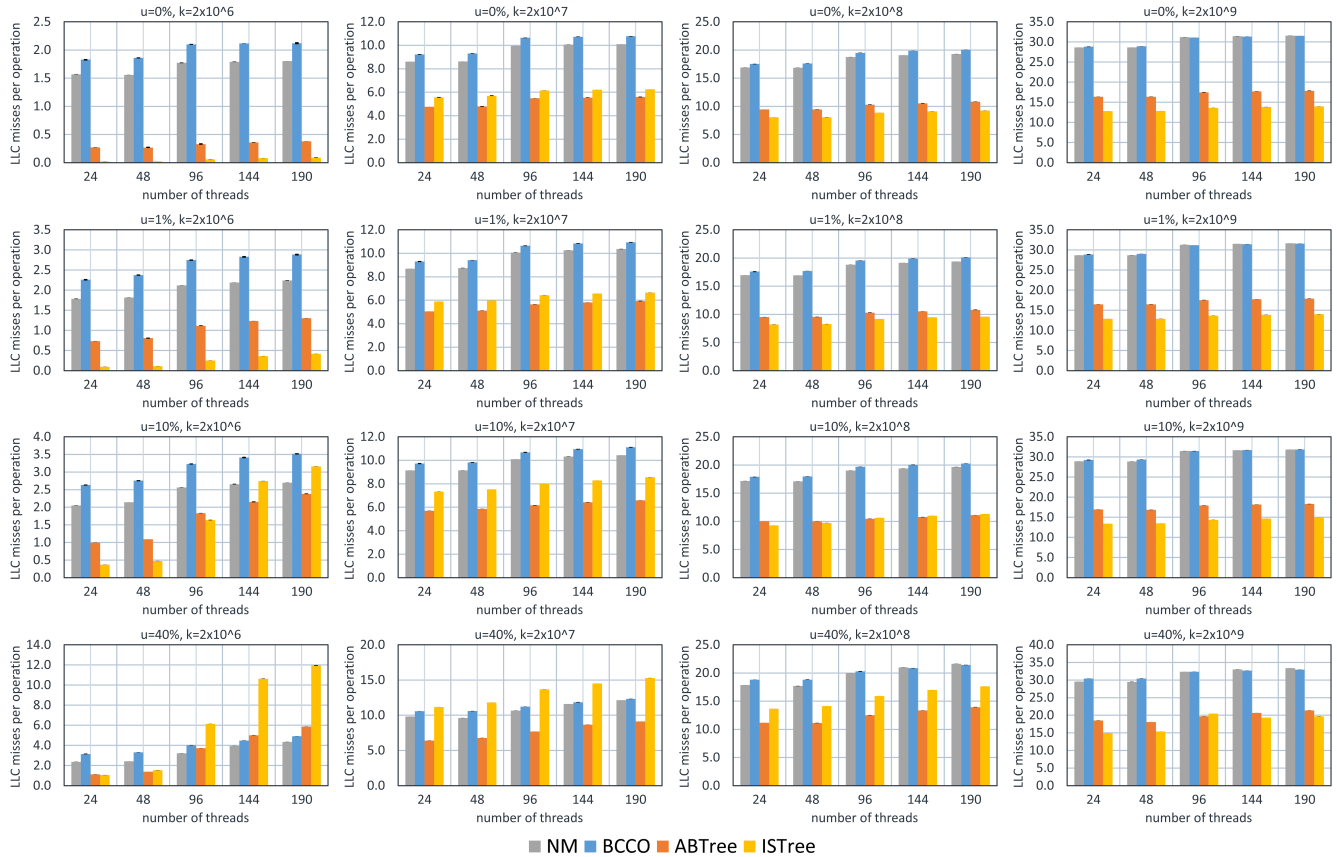
**Figure 11.** Summary of Last-Level Cache Misses, Lower is Better

## 6   Related Work

Sequential interpolation search was first proposed by Peterson [39], and subsequently analyzed by [24, 38, 58]. The *dynamic* case, where insertions and deletions are possible, was proposed by Frederickson [23]. The sequential IST variant we build on is by Mehlhorn and Tsakalidis [34]. This data structure supports amortized insertions and deletions in $O(\log n)$ time, under arbitrary distributions, and amortized insertion, deletion, and search, in $O(\log \log n)$ time under smoothness assumptions on the key distribution. To improve scalability, we augmented C-IST with parallel marking (to prevent updates during rebuilding), and a parallel rebuilding phase.

For *concurrent* search data structures ensuring predecessor queries, the work that is closest to ours is the SkipTrie [37], which allows predecessor queries in amortized expected $O(\log \log u + \gamma)$ steps, and insertions and deletions in $O(\gamma \log \log u)$ time, where $u$ is the size of the key space, and $\gamma$ is an upper bound on contention. The C-IST provides inferior runtime bounds in the worst case (e.g., $O(\log n)$ versus $O(\log \log u)$ amortized); however, the guarantees provided under distributional assumptions are asymptotically the same. We believe the C-IST should provide superior practical performance due to better cache behavior. We have

attempted to provide a comparison of the C-IST with an open-source implementation of the SkipTrie [1]; we found that this implementation had significant stability and performance issues, which render a fair comparison impossible.

There is considerable work on designing efficient concurrent search tree data structures with predecessor queries, e.g. [6, 9, 11, 13, 14, 14, 21, 36]. The average-case complexity of these operations is usually logarithmic in the number of keys. For large key counts (our target application) this search term dominates, giving the C-IST a significant performance advantage. This effect is apparent in our experimental section.

Other work on concurrent search tree data structures includes early work by Kung [31], Bronson's lock-based concurrent AVL trees [10], Pugh's concurrent skip list [54], and later improvements by Herlihy et al. [28] (which the JDK implementation is based on), non-blocking BSTs due to Ellen et al. [22], and the KiWi data structure due to Basin et al. [8].

The `DCSS` and `DCSS_READ` primitives that we rely on were originally proposed by Harris [27]. The `DCSS` primitive needs to allocate a descriptor object to synchronize multiple memory locations. Our C++ implementation of `DCSS`, due
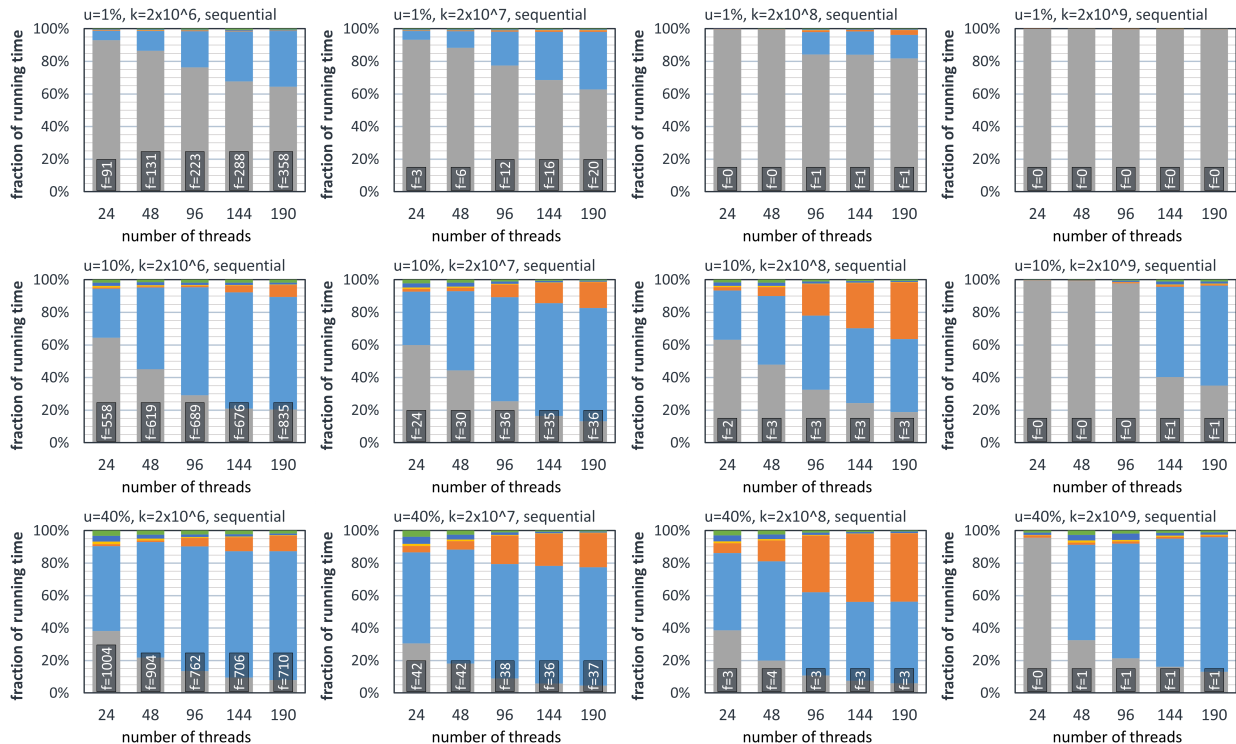
**Figure 12.** Breakdown of Time Spent, when Using Non-Collaborative Rebuilding (■ creating, ■ marking, ■ useless helping, ■ deallocation, ■ locating garbage, ■ other). Bars are annotated with $f$, the number of times the *root* (entire tree) was rebuilt.

to Arbel-Raviv and Brown [4], is able to recycle the descriptors. There are alternative primitives to `DCSS` with similar expressive power, such as the `GCAS` instruction [48], used to achieve snapshots in the Ctrie data structure.

Many concurrent data structures use the technique of snapshotting the entire data structure or some part thereof, with the goal of implementing a specific operation. The SnapQueue data structure [41] uses a *freezing* technique in which the writable locations are overwritten with special values such that the subsequent `CAS` operations fail. Ctries [40, 45, 46, 48] use the afore-mentioned `GCAS` operation to prevent further updates to the data structure. Work-stealing iterators [52, 53], used in work-stealing schedulers [33, 51] for data-parallel collections [47], use similar techniques to capture a snapshot of the iterator state.

The core motivation behind C-IST is to decrease the number of pointer hops during the key search. The underlying reason for this is that cache misses, which are incurred during the key search, are the dominating factor in the operation's running time. The motivation behind the recently proposed Cache-Trie data structure, a non-blocking Ctrie variant, is similar – Cache-Tries use an auxiliary, quiescently-consistent table to speed up the key searches [42–44].

Our implementation of the C-IST data structure uses a scalable concurrent counter in the root node to track the number of updates since the last rebuild of the root node.

In the past, a large body of research focused on scalable concurrent counters, both deterministic and probabilistic variant thereof [2, 7, 19, 20, 29, 57]. Scalable counters are useful in a number of other non-blocking data structures, which use counters to track their size or various statistics about the data structure. These include non-blocking queues [35], FlowPools [49, 50, 56], concurrent hash maps in the JDK [32], certain concurrent skip list implementations [25], and graphs with reachability queries [17].

Our C-IST implementation is done in C++, and it uses a custom concurrent memory management scheme due to Brown [12, 16]. In addition, our implementation uses techniques that decrease memory-allocator pressure by reusing the descriptors that are typically used in lock-free algorithms [3–5].

## 7  Conclusion

We proved the correctness, linearizability and non-blocking behavior of the C-IST data structure. We furthermore analyzed the doubly-logarithmic complexity of the basic C-IST operations. An extended experimental evaluation suggest that C-IST significantly improves upon the performance of competitive state-of-the-art search data structures, in some cases by up to ≈ 3.5×, and the current best-performing alternative by up to 50%. In highly skewed workloads, C-IST loses some of its performance advantages but nevertheless
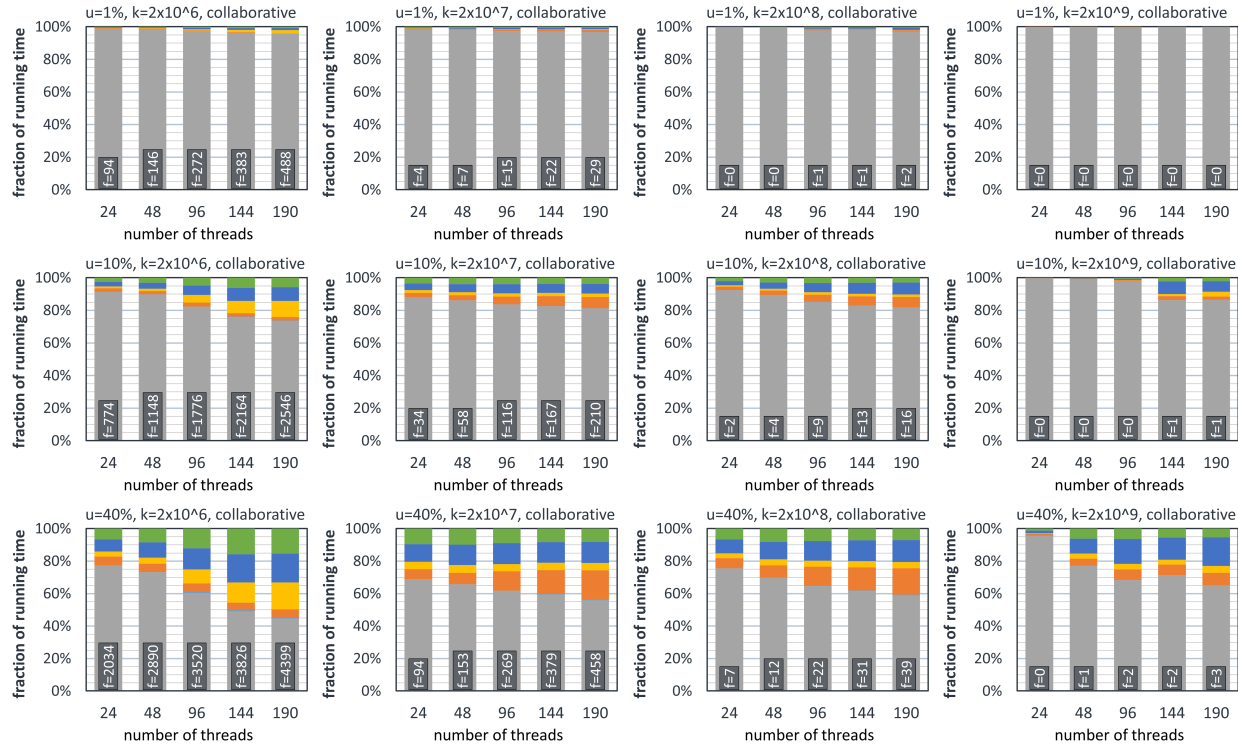
**Figure 13.** Breakdown of Time Spent, when Using Collaborative Rebuilding (■ creating, ■ marking, ■ useless helping, ■ deallocation, ■ locating garbage, ■ other). Bars are annotated with $f$, the number of times the *root* (entire tree) was rebuilt.
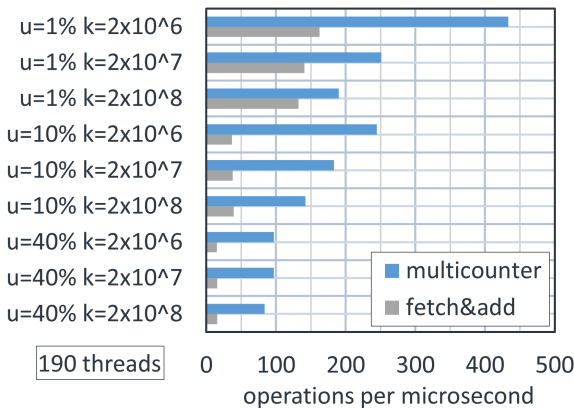


**Figure 14.** Effect of Using a Multicounter in the Root Node

exhibits a similar performance as the next best concurrent tree data structure.

Given the importance of Big Data workloads and very large databases, we believe that the techniques used in the C-IST data structure have the potential to trigger an intriguing line of future work.

## References

[1] 2018. SkipTrie Implementation at GitHub. https://github.com/JoeLeavitt/SkipTrie

[2] Dan Alistarh, Trevor Brown, Justin Kopinsky, Jerry Z. Li, and Giorgi Nadiradze. 2018. Distributively Linearizable Data Structures. In *Proceedings of the 30th on Symposium on Parallelism in Algorithms and Architectures (SPAA '18)*. ACM, New York, NY, USA, 133–142. https://doi.org/10.1145/3210377.3210411

[3] Maya Arbel-Raviv and Trevor Brown. 2017. POSTER: Reuse, don't Recycle: Transforming Algorithms that Throw Away Descriptors. In *Proceedings of the 22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, Austin, TX, USA, February 4-8, 2017*. 429–430. http://dl.acm.org/citation.cfm?id=3019035

[4] Maya Arbel-Raviv and Trevor Brown. 2017. Reuse, Don't Recycle: Transforming Lock-Free Algorithms That Throw Away Descriptors. In *31st International Symposium on Distributed Computing, DISC 2017, October 16-20, 2017, Vienna, Austria*. 4:1–4:16. https://doi.org/10.4230/LIPIcs.DISC.2017.4

[5] Maya Arbel-Raviv and Trevor Brown. 2017. Reuse, don't Recycle: Transforming Lock-free Algorithms that Throw Away Descriptors. *CoRR* abs/1708.01797 (2017). arXiv:1708.01797 http://arxiv.org/abs/1708.01797

[6] Maya Arbel-Raviv, Trevor Brown, and Adam Morrison. 2018. Getting to the Root of Concurrent Binary Search Tree Performance. In *USENIX Annual Technical Conference*.

[7] James Aspnes, Maurice Herlihy, and Nir Shavit. 1994. Counting Networks. *J. ACM* 41, 5 (Sept. 1994), 1020–1048. https://doi.org/10.1145/185675.185815

[8] Dmitry Basin, Edward Bortnikov, Anastasia Braginsky, Guy Golan-Gueta, Eshcar Hillel, Idit Keidar, and Moshe Sulamy. 2017. KiWi: A Key-Value Map for Scalable Real-Time Analytics. In *Proceedings of the 22Nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17)*. ACM, New York, NY, USA, 357–369. https://doi.org/10.1145/3018743.3018761
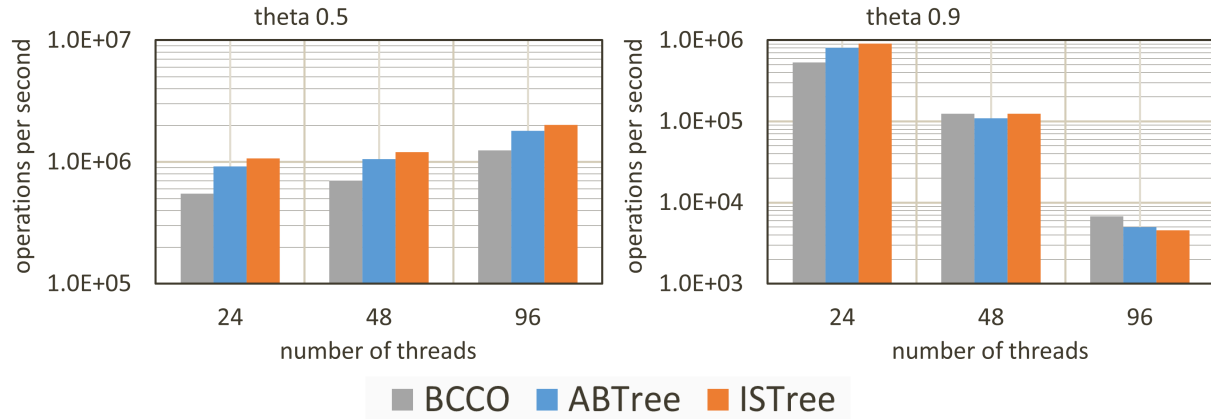
**Figure 15.** YCSB database performance with different index data structures, and a skewed *key access pattern*. (NM omitted because it is slower than BCCO.)
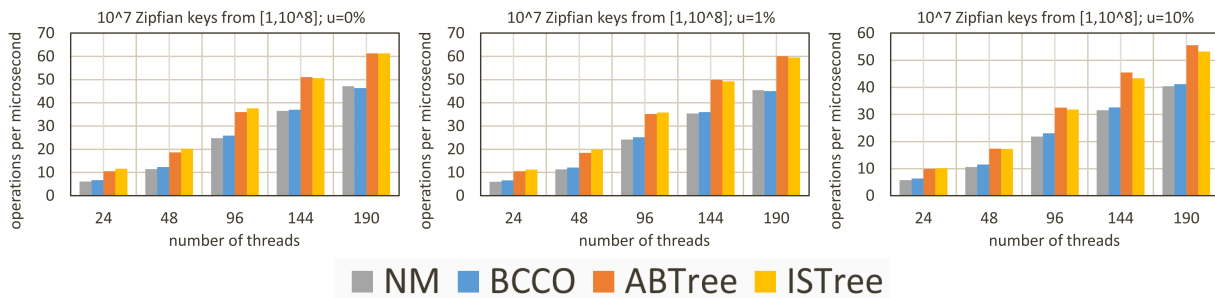


**Figure 16.** Synthetic benchmark in which the *set of keys stored* in a data structure is highly skewed.

[9] Anastasia Braginsky and Erez Petrank. 2012. A lock-free B+ tree. In *Proceedings of the twenty-fourth annual ACM symposium on Parallelism in algorithms and architectures*. ACM, 58–67.

[10] Nathan G. Bronson, Jared Casper, Hassan Chafi, and Kunle Olukotun. 2010. A Practical Concurrent Binary Search Tree. *SIGPLAN Not.* 45, 5 (Jan. 2010), 257–268. https://doi.org/10.1145/1837853.1693488

[11] Trevor Brown. 2014. B-slack Trees: Space Efficient B-Trees. In *Algorithm Theory - SWAT 2014 - 14th Scandinavian Symposium and Workshops, Copenhagen, Denmark, July 2-4, 2014. Proceedings*. 122–133. https://doi.org/10.1007/978-3-319-08404-6_11

[12] Trevor Brown. 2017. Reclaiming memory for lock-free data structures: there has to be a better way. *CoRR* abs/1712.01044 (2017). arXiv:1712.01044 http://arxiv.org/abs/1712.01044

[13] Trevor Brown. 2017. *Techniques for Constructing Efficient Data Structures*. Ph.D. Dissertation. University of Toronto.

[14] Trevor Brown and Hillel Avni. 2012. Range Queries in Non-blocking *k*-ary Search Trees. In *Principles of Distributed Systems, 16th International Conference, OPODIS 2012, Rome, Italy, December 18-20, 2012. Proceedings*. 31–45. https://doi.org/10.1007/978-3-642-35476-2_3

[15] Trevor Brown, Aleksandar Prokopec, and Dan Alistarh. 2020. Non-Blocking Interpolation Search Trees with Doubly-Logarithmic Running Time. In *Proceedings of the 25th Symposium on Principles and Practice of Parallel Programming (PPoPP '20)*. ACM, New York, NY, USA.

[16] Trevor Alexander Brown. 2015. Reclaiming Memory for Lock-Free Data Structures: There has to be a Better Way. In *Proceedings of the 2015 ACM Symposium on Principles of Distributed Computing, PODC 2015, Donostia-San Sebastián, Spain, July 21 - 23, 2015*. 261–270. https:

//doi.org/10.1145/2767386.2767436

[17] Bapi Chatterjee, Sathya Peri, Muktikanta Sa, and Nandini Singhal. 2019. A Simple and Practical Concurrent Non-Blocking Unbounded Graph with Linearizable Reachability Queries. In *Proceedings of the 20th International Conference on Distributed Computing and Networking (ICDCN '19)*. Association for Computing Machinery, New York, NY, USA, 168–177. https://doi.org/10.1145/3288599.3288617

[18] Tudor David, Rachid Guerraoui, and Vasileios Trigonakis. 2015. Asynchronized Concurrency: The Secret to Scaling Concurrent Search Data Structures. In *ASPLOS*.

[19] Damian Dechev and Bjarne Stroustrup. 2009. Scalable Nonblocking Concurrent Objects for Mission Critical Code. In *Proceedings of the 24th ACM SIGPLAN Conference Companion on Object Oriented Programming Systems Languages and Applications (OOPSLA '09)*. Association for Computing Machinery, New York, NY, USA, 597–612. https://doi.org/10.1145/1639950.1639954

[20] Dave Dice, Yossi Lev, and Mark Moir. 2013. Scalable Statistics Counters. In *Proceedings of the Twenty-Fifth Annual ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '13)*. Association for Computing Machinery, New York, NY, USA, 43–52. https://doi.org/10.1145/2486159.2486182

[21] Dana Drachsler, Martin Vechev, and Eran Yahav. 2014. Practical concurrent binary search trees via logical ordering. *ACM SIGPLAN Notices* 49, 8 (2014), 343–356.

[22] Faith Ellen, Panagiota Fatourou, Eric Ruppert, and Franck van Breugel. 2010. Non-blocking Binary Search Trees. In *Proceedings of the 29th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing (PODC '10)*. ACM, New York, NY, USA, 131–140. https:

//doi.org/10.1145/1835698.1835736

[23] Greg N Frederickson. 1983. Implicit data structures for the dictionary problem. *Journal of the ACM (JACM)* 30, 1 (1983), 80–94.

[24] Gaston H Gonnet, Lawrence D Rogers, and J Alan George. 1980. An algorithmic and complexity analysis of interpolation search. *Acta Informatica* 13, 1 (1980), 39–52.

[25] Vincent Gramoli. 2015. More than You Ever Wanted to Know about Synchronization: Synchrobench, Measuring the Impact of the Synchronization on Concurrent Algorithms. In *Proceedings of the 20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP 2015)*. Association for Computing Machinery, New York, NY, USA, 1–10. https://doi.org/10.1145/2688500.2688501

[26] Jim Gray, Prakash Sundaresan, Susanne Englert, Ken Baclawski, and Peter J. Weinberger. 1994. Quickly Generating Billion-record Synthetic Databases. In *Proceedings of the 1994 ACM SIGMOD International Conference on Management of Data (SIGMOD '94)*. ACM, New York, NY, USA, 243–252. https://doi.org/10.1145/191839.191886

[27] Timothy L. Harris, Keir Fraser, and Ian A. Pratt. 2002. A Practical Multi-word Compare-and-Swap Operation. In *Proceedings of the 16th International Conference on Distributed Computing (DISC '02)*. Springer-Verlag, Berlin, Heidelberg, 265–279. http://dl.acm.org/citation.cfm?id=645959.676137

[28] Maurice Herlihy, Yossi Lev, Victor Luchangco, and Nir Shavit. 2006. A Provably Correct Scalable Concurrent Skip List.

[29] Maurice Herlihy, Beng-Hong Lim, and Nir Shavit. 1995. Scalable Concurrent Counting. *ACM Trans. Comput. Syst.* 13, 4 (Nov. 1995), 343–364. https://doi.org/10.1145/210223.210225

[30] Shane V. Howley and Jeremy Jones. 2012. A Non-blocking Internal Binary Search Tree. In *SPAA*.

[31] H. T. Kung and Philip L. Lehman. 1980. Concurrent Manipulation of Binary Search Trees. *ACM Trans. Database Syst.* 5, 3 (Sept. 1980), 354–382. https://doi.org/10.1145/320613.320619

[32] Doug Lea. 2018. Doug Lea's Workstation. http://g.oswego.edu/

[33] Jonathan Lifflander, Sriram Krishnamoorthy, and Laxmikant V. Kale. 2013. Steal Tree: Low-Overhead Tracing of Work Stealing Schedulers. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*. Association for Computing Machinery, New York, NY, USA, 507–518. https://doi.org/10.1145/2491956.2462193

[34] Kurt Mehlhorn and Athanasios Tsakalidis. 1993. Dynamic interpolation search. *Journal of the ACM (JACM)* 40, 3 (1993), 621–634.

[35] Maged M. Michael and Michael L. Scott. 1996. Simple, Fast, and Practical Non-Blocking and Blocking Concurrent Queue Algorithms. In *Proceedings of the Fifteenth Annual ACM Symposium on Principles of Distributed Computing (PODC '96)*. Association for Computing Machinery, New York, NY, USA, 267–275. https://doi.org/10.1145/248052.248106

[36] Aravind Natarajan and Neeraj Mittal. 2014. Fast Concurrent Lock-free Binary Search Trees. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '14)*. 317–328.

[37] Rotem Oshman and Nir Shavit. 2013. The SkipTrie: Low-depth Concurrent Search Without Rebalancing. In *Proceedings of the 2013 ACM Symposium on Principles of Distributed Computing (PODC '13)*. ACM, New York, NY, USA, 23–32. https://doi.org/10.1145/2484239.2484270

[38] Yehoshua Perl and Edward M Reingold. 1977. Understanding the complexity of interpolation search. *Inform. Process. Lett.* 6, 6 (1977), 219–222.

[39] W Wesley Peterson. 1957. Addressing for random-access storage. *IBM journal of Research and Development* 1, 2 (1957), 130–146.

[40] Aleksandar Prokopec. 2014. Data Structures and Algorithms for Data-Parallel Computing in a Managed Runtime. (2014).

[41] Aleksandar Prokopec. 2015. SnapQueue: Lock-free Queue with Constant Time Snapshots. In *Proceedings of the 6th ACM SIGPLAN Symposium on Scala (SCALA 2015)*. ACM, New York, NY, USA, 1–12.

https://doi.org/10.1145/2774975.2774976

[42] Aleksandar Prokopec. 2017. Analysis of Concurrent Lock-Free Hash Tries with Constant-Time Operations. *ArXiv e-prints* (Dec. 2017). arXiv:cs.DS/1712.09636

[43] Aleksandar Prokopec. 2018. Cache-tries: Concurrent Lock-free Hash Tries with Constant-time Operations. In *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '18)*. ACM, New York, NY, USA, 137–151. https://doi.org/10.1145/3178487.3178498

[44] Aleksandar Prokopec. 2018. *Efficient Lock-Free Removing and Compaction for the Cache-Trie Data Structure*. Springer International Publishing, Cham.

[45] Aleksandar Prokopec, Phil Bagwell, and Martin Odersky. 2011. *Cache-Aware Lock-Free Concurrent Hash Tries*. Technical Report.

[46] Aleksandar Prokopec, Phil Bagwell, and Martin Odersky. 2011. *Lock-Free Resizeable Concurrent Tries*. Springer Berlin Heidelberg, Berlin, Heidelberg, 156–170. https://doi.org/10.1007/978-3-642-36036-7_11

[47] Aleksandar Prokopec, Phil Bagwell, Tiark Rompf, and Martin Odersky. 2011. A Generic Parallel Collection Framework. In *Proceedings of the 17th international conference on Parallel processing - Volume Part II (Euro-Par'11)*. Springer-Verlag, Berlin, Heidelberg, 136–147. http://dl.acm.org/citation.cfm?id=2033408.2033425

[48] Aleksandar Prokopec, Nathan Grasso Bronson, Phil Bagwell, and Martin Odersky. 2012. Concurrent Tries with Efficient Non-blocking Snapshots. In *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '12)*. ACM, New York, NY, USA, 151–160. https://doi.org/10.1145/2145816.2145836

[49] Aleksandar Prokopec, Heather Miller, Philipp Haller, Tobias Schlatter, and Martin Odersky. 2012. *FlowPools: A Lock-Free Deterministic Concurrent Dataflow Abstraction, Proofs*. Technical Report.

[50] Aleksandar Prokopec, Heather Miller, Tobias Schlatter, Philipp Haller, and Martin Odersky. 2012. FlowPools: A Lock-Free Deterministic Concurrent Dataflow Abstraction. In *LCPC*. 158–173.

[51] Aleksandar Prokopec and Martin Odersky. 2014. *Near Optimal Work-Stealing Tree Scheduler for Highly Irregular Data-Parallel Workloads*. Springer International Publishing, Cham, 55–86. https://doi.org/10.1007/978-3-319-09967-5_4

[52] Aleksandar Prokopec, Dmitry Petrashko, and Martin Odersky. 2014. On Lock-Free Work-stealing Iterators for Parallel Data Structures. (2014), 10. http://infoscience.epfl.ch/record/196627

[53] A. Prokopec, D. Petrashko, and M. Odersky. 2015. Efficient Lock-Free Work-Stealing Iterators for Data-Parallel Collections. In *2015 23rd Euromicro International Conference on Parallel, Distributed, and Network-Based Processing*. 248–252. https://doi.org/10.1109/PDP.2015.65

[54] William Pugh. 1990. *Concurrent Maintenance of Skip Lists*. Technical Report. College Park, MD, USA.

[55] Arunmoezhi Ramachandran and Neeraj Mittal. 2015. A Fast Lock-Free Internal Binary Search Tree. In *ICDCN*.

[56] Tobias Schlatter, Aleksandar Prokopec, Heather Miller, Philipp Haller, and Martin Odersky. 2012. Multi-Lane FlowPools: A Detailed Look. (2012), 13.

[57] Guy L. Steele and Jean-Baptiste Tristan. 2016. Adding Approximate Counters. In *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '16)*. Association for Computing Machinery, New York, NY, USA, Article Article 15, 12 pages. https://doi.org/10.1145/2851141.2851147

[58] Andrew C Yao and F Frances Yao. 1976. The complexity of searching an ordered random table. In *Foundations of Computer Science, 1976., 17th Annual Symposium on*. IEEE, 173–177.

[59] Xiangyao Yu, George Bezerra, Andrew Pavlo, Srinivas Devadas, and Michael Stonebraker. 2014. Staring into the Abyss: An Evaluation of Concurrency Control with One Thousand Cores. *VLDB* 8, 3 (Nov. 2014).