

MIT Open Access Articles

Leaplist: lessons learned in designing tm-supported range queries

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Hillel Avni, Nir Shavit, and Adi Suissa. 2013. Leaplist: lessons learned in designing tm-supported range queries. In Proceedings of the 2013 ACM symposium on Principles of distributed computing (PODC '13). ACM, New York, NY, USA, 299-308.

As Published: <http://dx.doi.org/10.1145/2484239.2484254>

Publisher: Association for Computing Machinery (ACM)

Persistent URL: <http://hdl.handle.net/1721.1/90890>

Version: Original manuscript: author's manuscript prior to formal peer review

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



Leaplist: Lessons Learned in Designing TM-Supported Range Queries

Hillel Avni

Tel-Aviv University
Tel-Aviv 69978, Israel
hillel.avni@gmail.com

Nir Shavit

MIT and Tel-Aviv University
shanir@csail.mit.edu

Adi Suissa

Department of Computer Science
Ben-Gurion University of the Negev
Be'er Sheva, Israel
adisuis@cs.bgu.ac.il

Abstract

We introduce *Leap-List*, a concurrent data-structure that is tailored to provide linearizable range queries. A lookup in *Leap-List* takes $O(\log n)$ and is comparable to a balanced binary search tree or to a skip-list. However, in *Leap-List*, each node holds up-to K immutable key-value pairs, so collecting a linearizable range is K times faster than the same operation performed non-linearizably on a skip-list.

We show how software transactional memory support in a commercial compiler helped us create an efficient lock-based implementation of *Leap-List*. We used this STM to implement short transactions which we call Locking Transactions (LT), to acquire locks, while verifying that the state of the data-structure is legal, and combine them with a transactional COP mechanism to enhance data structure traversals.

We compare *Leap-List* to prior implementations of skip-lists, and show that while updates in the *Leap-List* are slower, lookups are somewhat faster, and for range-queries the *Leap-List* outperforms the skip-list's non-linearizable range query operations by an order of magnitude. We believe that this data structure and its performance would have been impossible to obtain without the STM support.

Keywords Transactional-Memory, Data-Structures, Range-Queries

1. Introduction and Related Work

We are interested in linearizable concurrent implementations of an abstract dictionary data structure that stores key-value pairs and supports, in addition to the usual $\text{Update}(\text{key}, \text{value})$, $\text{Remove}(\text{key})$, and $\text{Find}(\text{key})$, a $\text{Range-Query}(a, b)$ operation, where $a \leq b$, which returns all pairs with keys in the closed interval $[a, b]$, where a and b may not be in the data structure. This type of data structure is useful for various database applications, in particular in-memory databases. This paper is interested in the design of high performance linearizable concurrent range queries. As such, the typically logarithmic search for the first item in the range is not the most important performance element. Rather, it is the coordination and synchronization around the sets of neighboring keys being collected in the sequence. This is a tricky new synchronization problem and our goal is to evaluate which transactional support paradigm, if any, can help in attaining improved performance for range queries.

1.1 Related Work

Perhaps the most straightforward way to implement a linearizable concurrent version of an abstract dictionary-with-range-queries, is to directly employ software transactional memory (STM) in implementing its methods. An STM allows a programmer to specify

that certain blocks of code should be executed atomically relative to one another. Recently, several fast concurrent binary search-tree algorithms using STM have been introduced by Afek et al. [2] and Bronson et al. [4]. Although they offer good performance for Updates, Removes and Finds, they achieve this performance, in part, by carefully limiting the amount of data protected by the transactions. However, as we show empirically in this paper, computing a range query means protecting all keys in the range from modification during a transaction, leading to poor performance using the direct STM approach.

Another simple approach is to lock the entire data structure, and compute a range query while it is locked. One can refine this technique by using a more fine-grained locking scheme, so that only part of the data structure needs to be locked to perform an update or compute a range query. For instance, in leaf-oriented trees, where all key-value pairs in the set are stored in the leaves of the tree, updates to the tree can be performed by local modifications close to the leaves. Therefore, it is often sufficient to lock only the last couple of nodes on the path to a leaf, rather than the entire path from the root. However, as was the case for STM, a range query can only be computed if every key in the range is protected, so typically every node containing a key in the range must be locked.

Brown and Avni [5] introduced range queries in k -ary trees with immutable key. The k -ary trees allow efficient range-queries by collecting nodes in a depth-first-search order, followed by a validation stage. The nodes are scanned, and if any one is outdated, the process is retried from the start. The k -ary search tree is not balanced, and its operations cannot be composed.

Ctrie is a non-blocking concurrent hash trie, which offers $O(1)$ time snapshot, due to Prokopec et al. [10]. Keys are hashed, and the bits of these hashes are used to navigate the trie. To facilitate the computation of fast snapshots, a sequence number is associated with each node in the data structure. Each time a snapshot is taken, the root is copied and its sequence number is incremented. An update or search in the trie reads this sequence number seq when it starts and, while traversing the trie, it duplicates each node whose sequence number is less than seq . The update then performs a variant of a double-compare-single-swap operation to atomically change a pointer while ensuring the roots current sequence number matches seq . Because keys are ordered by their hashes in the trie, it is hard to use Ctrie to efficiently implement range queries. To do so, one must iterate over all keys in the snapshot.

The B-Tree data structure can be used for range queries, however, when looking at the concurrent versions of B-Trees such as the lock-free one of Braginsky and Petrank [3], and the blocking, industry standard from [12], both do not support the range-query functionality. Both algorithms do not have leaf-chaining, forcing one to perform a sequence of lookups to collect the desired range.

In [12] this would imply holding a lock on the root for a long time, and in [3] it seems difficult to get a linearizable result. In addition, the keys in both are mutable so one would have to copy each entry individually.

1.2 The Leap-List in a Nutshell

Leap-Lists are Skip-Lists [11] with “fat” nodes and an added shortcut access mechanism in the style of the String B-tree of Ferragina and Grossi [6]. They have the same probabilistic guarantee for balancing, and the same layered forward pointers as Skip-Lists. Each *Leap-List* node holds up to K immutable keys from a specific range, and an immutable bitwise trie is embedded in each node to facilitate fast lookups when K is large.

When considering large range queries, the logarithmic-time lookup for the start of the range accounts for only a small part of the operation’s complexity. Especially when the whole structure resides in memory. The design complexity of a full k -ary structure (in which nodes at all levels have K elements), with $\log_k(n)$ lookup time is thus not justified. In our *Leap-List*, unlike full k -ary structures, an update implies at most one split or merge of a node, and only at the leaf level. This allows updates to lock only the specific leaf being changed and only for the duration of changing pointers from the old node to the new one.

For *Leap-List* synchronization, we checked the following options, sorted in an increasing order of required effort:

- **Pure STM:** We tried to put each *Leap-List* operation in a software transactional memory (STM) transaction. This option was especially attractive with the rising support for STM in mainstream compilers. Unfortunately, as we report, we discovered that this approach introduced unacceptable overheads.
- **Read-write locks:** We explored read-write locks per *Leap-List*. While the read-locks were very scalable, the write locks serialized many workloads, hence making updates relatively slow.
- **COP:** We employed consistency oblivious programming (COP) [2] to reduce the overhead of STM. In COP, the read-only prefix of the operation is executed without any synchronization, followed by an STM transaction that checks the correctness of the prefix execution and performs any necessary updates. The COP requires that an un-instrumented traversal of the structure will not crash, which implies strong isolation of transactions in the underlying STM. Otherwise the traversal encounters uncommitted data, and hitting uncommitted data inevitably leads to uninitialized pointers, unallocated buffers, and segmentation faults. The current GCC-TM compiler uses weakly isolated transactions. Thus, we had to add transactions also in read-only regions of the code which hurt performance.
- **Locking Transactions (LT):** With LT, transactions are used only to acquire locks, and not to write tentative data. Thus, a read which sees unlocked data knows it is committed. Another aspect of LT, is that using a short transaction anyone can lock any data and use it.

We use LT to improve the performance of the previous COP algorithm. In the COP, an updating operation performs its read-only prefix without synchronization, and then executes the verification and updates inside a transaction. In LT, the read-only part is checking for locks, and retries. These checks have negligible overhead compared to a transaction. Then the transaction atomically verifies validity and locks the written addresses. After the transaction commits, a postfix of the operation writes the data to the locked locations and releases them.

- **Fine grained locks:** To generate the fine-grain version of LT *Leap-List* we had to recreate mechanisms that exist in STM, and still, did not manage to create a fully stable implementation.

In case of a merge, where a remove replaces two old nodes by one new node, we need to lock all pointers to and from both nodes. Here, unlike the skip-list case [9], locking can fail at any point and force us to release all locks and retry to avoid deadlocks. This unrolling is “free” using an STM.

Once a set of nodes is locked, a thread needs to perform validations on the state of the data structure, such as checking live marks etc. With LT, using STM, these validations happen before acquiring the locks, and then when committing, an abort will happen if any check should fail. Thus the locks are taken for a shorter duration. To improve our performance we would need to execute a form of STM revalidation.

After executing the above sequence, we found that our fine grained implementation still suffered from live-locks; we did not manage to avoid them. These live-locks were eliminated with the STM based LT approach.

Our conclusion was that we were effectively reproducing the very mechanisms that are already given by an STM, and still did not get the stability of an STM. The LT *Leap-List* implementation has minimal overhead because lookups do not execute transactions and range-queries execute one instrumented access per K values in the range. The LT *Leap-List* is thus the most effective solution.

This paper is organized as follows. Section 2 gives a detailed description of the *Leap-List* design and operations’ implementation. In Section 3 we show the LT technique is the best performer for *Leap-List* synchronization, and is scalable even when transactions encompass operations on multiple *Leap-Lists*. Finally, in Section 4 we summarize our work, and give some directions for future work.

2. Leap-List Design

We now describe the detailed design of our *L-Leap-Lists* data structure. Note that the updating functions compose operations on multiple *Leap-Lists*. Our implementation supports the following operations:

- **Update(ll, k, v, s)** - Receives arrays of *Leap-Lists*, keys and values of size s , and updates the value of the key $k[i]$ to be $v[i]$ in *Leap-List* $ll[i]$. If the key $k[i]$ is not present in $ll[i]$, it is added to $ll[i]$ with the given value $v[i]$.
- **Remove(ll, k, s)** - Receives arrays of *Leap-Lists* and keys of size s , and removes the key-value pair of the given key $k[i]$ from $ll[i]$.
- **Lookup(l, k)** - Receives a single *Leap-List* and a key, and returns the value of the corresponding given key k in l . The operation returns an indication in case the key is not present in l .
- **Range-Query(l, k_{from}, k_{to})** - Receives a single *Leap-List* and 2 keys, and returns the values of all keys in l which are in the range $[k_{from}, k_{to}]$.

The *Update* and *Remove* are linearizable operations applied to *L-Leap-Lists* and *Lookup* and *Range-Query* search are linearizable operations applied to a single *Leap-List*. This allows concurrent operations on multiple database table indexes. We implemented the *L-Leap-List* data-structure using the experimental GCC Transactional Memory (GCC-TM). GCC-TM is a word-based STM implementation with a default configuration in which transactions are weakly isolated.

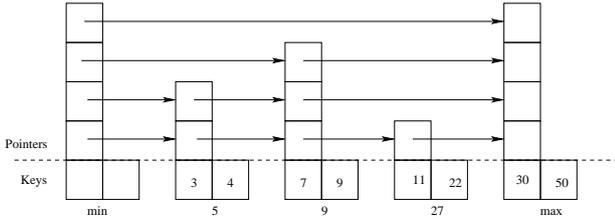


Figure 1: A single *Leap-List* with maximum height of 4 and node size of 2. The number below each node is the highest possible key of that node. The left-most node is always empty.

Leap-List Node data-structure

```

1 define Node: struct {
2   live: boolean,
3   high: unsigned long,
4   count: unsigned long,
5   level: byte,
6   next: array of *node,
7   values: {key-value} sorted_pairs,
8   trie: {key-index in node} trie
9 };

```

Figure 2: Data-Structure description

An example of a single *Leap-List* structure is depicted in Figure 1. As shown, each node may hold up to a predefined number of data items with different keys and a set of pointers for each level below that node level. The operations of the *Leap-Lists* are implemented using the COP scheme [2], where in the “search phase”¹ the data-structure pointers are accessed outside of a transactional context to achieve better performance. Due to the weakly isolated nature of GCC-TM and the need to prevent uncommitted pointers from guiding us to uncommitted nodes, we use a novel method of writing marked pointers in a transaction and removing the mark after a successful commit. A marked pointer indicates that the pointed node is currently being updated by an active transaction, or was updated by a transaction that was successfully committed. Another alternative we explored was to access pointers in single-location read transactions. However, this alternative proved to have a larger negative impact on performance with the current GCC-TM implementation. Nevertheless, we expect it will exhibit the best performance with HTM support, as single-location uncontended read transactions should be ideal for HTM.

2.1 Leap-List Data-Structure

The *Leap-List* node is presented in Figure 2. It holds a *live* mark, for COP verification; *high*, which denotes the upper bound of its keys range; *count*, which is the number of key-value pairs present in the node, and *level* which is the same as a level in skip-list. It also holds an array of forward pointers *next* each pointing to the next element in the corresponding level. A *trie* is used to quickly find the index of key *k* in the keys-values array, a technique introduced in the String B-tree of Ferragina and Grossi [6]. Note that unlike in a skip-list, where each node represents a single key, in *Leap-List* each node represents a range of keys, i.e. all the keys from a certain range. The *keys-values* array of size *count* holds all the keys and their corresponding values in the node. The *trie* uses the minimal number of levels to represent all the keys in the node, where the

¹We consider the “search phase” to consist of the Lookup and Range-Query operations, and the read-only accesses that are done in the Update and Remove operations before any write access.

Leap-List Search Predecessors

```

input : LeapList l, key k
output: Two node arrays of pointers - pa and na
10 node *x, *x_next;
11 int i;
12 retry:
13 x := l;
14 for i = max_level- 1; i ≥ 0; i = i - 1 do
15   while (true) do
16     x_next := x→next[i];
17     if MARKED(x_next) ∨ (¬x_next→live) then goto
    retry;
18     if x_next→high ≥ k then
19       break;
20     else
21       x := x_next
22   end
23 end
24 pa[i] := x;
25 na[i] := x_next;
26 end
27 return (pa, na);

```

Figure 3: *Leap-List* Search Predecessors operation

lowest level is comprised of indexes of the keys’ values in the *keys-values* array.

Upon initialization the empty list contains two nodes. One, sentinel node, whose range is bounded from above by $-\infty$, and the second, which has no keys and a high value of ∞ (and thus encompasses the range $(-\infty, \infty)$). The level of the sentinel node is the maximal level (**max_level**), and its next pointers all point to the second node. Two consecutive nodes, define the range encompassed by the second node. If node N_2 follows N_1 , then the range of N_2 is $(N_1.high \text{ to } N_2.high]$.

In *Leap-List*, a node’s keys-values array is immutable, and never changes after an update. We do this to support consistent range-query operations. When the key or value (and possibly the encompassed range) of a node is updated (due to an update or a remove operation), that node is replaced by a newer node with the modified keys-values array. If the node is full (i.e., the number of keys in the node reaches some predefined number), it is split into two consecutive nodes and the upper bound of the lower node is determined by the highest value in it. In case the modified node and its subsequent node are sparse (the number of keys in both nodes is less than some predefined number), the nodes are merged into a single node.

Unlike to a key lookup operation which returns a single value, a range-query operation returns an array of nodes, which are part of a consistent snapshot and hold all the values in the given range. In the rest of this section we will describe the *Leap-List* functions.

2.1.1 Searching for Predecessors

The search predecessors function from Figure 3 receives a key *k*, and traverses the *Leap-List* until the node *N* (that encompasses the range where key *k* is included) is reached. The function returns two arrays of nodes each, *pa* and *na*, each of size **max_level**. The *pa* array includes all the nodes that “immediately precede” node *N*. That is for each level *i* up to *N*’s level, $pa[i] \rightarrow next[i]$ points to *N*, and for levels higher than *N*’s level, the nodes that encompass keys that are smaller than *k* and their next pointer at level *i* points to a node with higher keys than *k*. The *na* array includes all the nodes that are adjacent to the *pa* nodes, and encompass keys that are greater-than-or-equal-to *k* (thus $na[i] \rightarrow next[i]$ is *N* for all

Leap-List Lookup

```
input : LeapList l, key k
output: Value or  $\perp$ 
28 node *na[max_level ];
29 (null, na) $\leftarrow$ PredecessorsSearch(l,k);
30 return (na[0] $\rightarrow$ values[get_index(na[0] $\rightarrow$ trie,k)].value);
```

Figure 4: Leap-List Lookup operation

levels up to N 's level). This function is used in the lookup and range-query operations, as well as in the beginning of the update and remove operations.

The traversal only compares the *high* key of the node in line 18 and decides if it should continue or stop at that node. When reading a pointer, the thread verifies that that pointer is not marked and that the node is still live in line 17, so it only traverses committed and valid nodes. (As previously noted, an alternative method would be replacing the mark by executing line 16 in a transaction. However, with the current GCC-TM implementation the overhead of starting a transaction is too high. We estimate that with HTM this would work much better, and will actually make the lookup wait-free, as a single-location read transaction must succeed.)

2.1.2 Lookups

The lookup operation is presented in Figure 4, and is using the predecessors search function. Note that the node returned in $na[0]$ is the node that has k in its range. We can prove the lookup is linearizable, as the predecessors search traverses only committed nodes. If a thread searches for the key k , it must traverse a node that k is in its range, and if such a live node is reached, then this node was present in the data-structure during the lookup execution.

In line 30, *Lookup* uses the node's trie to extract the index of the value of key k in the array values, and returns the value from that index.

2.1.3 Range Queries

The range query operation is presented in Figure 5, and starts with a predecessors search to find the node where the range starts from. Then, within a transaction, it first checks that the node is still live in line 39 and if not aborts, and retries the range-query operation in line 45. If the node is still marked as live, the transaction traverses the lowest level of the *Leap-List*'s pointers from the first node to the node which has a *high* value which is higher than the requested range high bound, and retrieves a snapshot range query. Note that in line 41 the algorithm ensures that even in the case of a partial update to the pointer to the next node (due to update or remove operations), it can still traverse through it.

2.1.4 Updates

Figure 6 describes the update function. As previously described, the function receives arrays of *Leap-Lists*, keys and values, and their size. The update operation either inserts a new key-value pair to each *Leap-List* if the key is not already present, or otherwise updates the key's value.

The function is divided into the following 3 parts: (1) setup (Figure 8), (2) LT (Figure 9), and (3) release and update (Figure 10). During the setup part, a thread iterates over each *Leap-List*, performs a predecessors search, and creates a new node with its key-value pairs (including the updated key-value pair). Note that in case the number of keys in the node is above some threshold, it *splits* that node. During a split it creates 2 nodes: one with a new random height that holds the first half of the key-value pairs, and another with the same height as the old node that holds the second half of the key-value pairs. The *max_level* is set to the maximum

Leap-List Range Query

```
input : Leap list l, key low, key high
output: Set S of nodes
31 node *na[max_level ], *n;
32 boolean committed  $\leftarrow$  false;
33 retry:
34 S $\leftarrow$   $\emptyset$ ;
35 (null, na) $\leftarrow$ Search(l,low);
36 n := na[0];
37 tx_start;
38 while n $\rightarrow$ high<high do
39   if  $\neg$ n $\rightarrow$ live then tx_abort;
40   add(S,na[0]);
41   n := unmark((n $\rightarrow$ next[0]));
42 end
43 committed := true;
44 tx_end;
45 if  $\neg$ committed then goto retry;
46 return S;
```

Figure 5: Leap-List Range Query.

between the heights of the two nodes. The *CreateNewNodes* function updates the new node (nodes) with its (their) key-value pairs.

The LT part is executed in a single transaction. The algorithm again iterates over each *Leap-List* and first verifies that the updated node is still live (line 95), that all the predecessors' next pointers point to that node, and that the next pointers from that node are still valid (lines 96-104). (In case of a split, the algorithm also verifies this up to the *max_level* height.) In lines 105-111 it continues to verify and mark the pointers to the node and from the node, and in case of a split the nodes to and from the nodes up to the *max_height*. Note that if one of the conditions does not hold, the transaction is aborted, and the whole operation restarts. It finishes the transaction by setting the old node's live bit to false (line 113), and attempting to commit the transaction. We note that in this part, the transaction does not observe partial modifications made by other transactions, and so a successful commit ensures a consistent view of the nodes that are affected by the operation.

Following a successful transaction commit, the third part releases and updates the pointers of the predecessor nodes to point to the new node (nodes). In lines 116-137 the algorithm sets the next pointer of the new node (nodes) to the previous nodes that were in the *Leap-List*. It continues by setting the next pointers of the predecessor nodes to the new node (nodes) in lines 139-145, and finishes by setting the live flags of the new node (nodes) to true.

2.1.5 Remove

The remove function is presented in Figure 7. The function receives arrays of *Leap-Lists*, keys and their size, and linearizably removes the key-value pair of each given key from its corresponding *Leap-List*. In case a key is not found in a *Leap-List*, that *Leap-List* is not modified.

Similarly to the update function, the remove function is also divided to the setup (Figure 11), LT (Figure 12) and release and update (Figure 13) parts. During the setup part, the thread again iterates over each *Leap-List*, performs a predecessors search, and searches for the key to be removed. If a *Leap-List* does not contain the corresponding key, it moves on to the next *Leap-List*. In case the key exists it keeps the node that holds the key and its successor node in the *old_node* variables (line 154-161). The node and its adjacent node are merged if the sum of the key-value pairs in both nodes is below some threshold. It then verifies that the node and the

Leap-Lists Update

```

input : Leap-Lists ll, keys k, values v, and size s
47 node *pa[max_lists][max_level], *na[max_lists]
   [max_level], *n[max_lists];
48 node *new_node[max_lists][2];
49 int max_height[max_lists];
50 boolean committed := false, split[max_lists];
51 foreach j < s do
52   new_node[j][0] := new node;
53   new_node[j][1] := new node;
54 end
55 retry:
56 Update_Setup(ll, k, v, s, pa, na, n, new_node, max_height,
   split);
57 tx_start ;
58 Update_LT(s, pa, na, n, new_node, max_height);
59 committed := true;
60 tx_end;
61 if  $\neg$ committed then goto retry;
62 Update_Release_and_Update(s, pa, na, n, new_node, split);
63 Deallocate_unneeded_nodes.

```

Figure 6: Leap-List Update

Leap-List Remove

```

input : Leap-Lists ll, keys k, size s
64 node *pa[max_lists][max_level], *na[max_lists]
   [max_level], *n[max_lists];
65 node *old_node[max_lists][2];
66 boolean committed := false, merge[max_lists],
   changed[max_lists];
67 foreach j < s do
68   n[j] := new node;
69 end
70 retry_all:
71 Remove_Setup(ll, k, v, s, pa, na, n, old_node, merge,
   changed);
72 tx_start;
73 Remove_LT(s, pa, na, n, old_node, merge, changed);
74 committed := true;
75 tx_end ;
76 if  $\neg$ committed then goto retry_all;
77 Remove_Release_and_Update(s, pa, na, n, old_node, merge,
   changed);
78 Deallocate_unneeded_nodes.

```

Figure 7: Leap-List Remove

Leap-List Update - Setup

```

input : Leap-Lists ll, keys k, values v, size s, nodes pa,
   nodes na, nodes n, nodes new_node, integers
   max_height, booleans split
79 foreach j < s do
80   (pa[j],na[j]) ← PredecessorSearch(ll[j],k[j]);
81   n[j] := na[j][0];
82   if n[j] → count = node_size then
83     split[j] := true;
84     new_node[j][1] → level := n[j] → level;
85     new_node[j][0] → level := get_level();
86     max_height[j] := max(new_node[j][0] → level,
   new_node[j][1] → level);
87   else
88     split[j] := false;
89     new_node[j][0] → level := n[j] → level;
90     max_height[j] := new_node[j][0] → level;
91   end
92   CreateNewNodes(new_node[j], n[j], k[j], v[j], split[j]);
93 end

```

Figure 8: Leap-List Update - Setup.

Leap-List Update - LT

```

input : size s, nodes pa, nodes na, nodes n, nodes new_node,
   integers max_height
94 foreach j < s do
95   if  $\neg$ n[j] → live then tx_abort;
96   foreach i < n[j] → level do
97     if pa[j][i] → next[i] ≠ n[j] then tx_abort;
98     if  $\neg$ n[j] → next[i] → live then tx_abort;
99   end
100  foreach i < max_height[j] do
101    if pa[j][i] → next[j][i] ≠ na[j][i] then tx_abort;
102    if  $\neg$ pa[j][i] → live then tx_abort;
103    if  $\neg$ na[j][i] → live then tx_abort;
104  end
105  foreach i < n[j] → level do
106    if MARKED(n[j] → next[i]) then tx_abort;
107    n[j] → next[i] := MARK(n[j] → next[i]);
108  end
109  foreach i < max_height[j] do
110    if MARKED(pa[j][i] → next[i]) then tx_abort;
111    pa[j][i] → next[i] := MARK(pa[j][i] → next[i]);
112  end
113  n[j] → live := false ;
114 end

```

Figure 9: Leap-List Update - LT.

adjacent node (upon merge) are live, and if not, the retry of the last key removal from the current Leap-List is performed. The thread concludes this part by calling *RemoveAndMerge* which updates a new node with the key-value pairs from the node (and the adjacent node), without the removed key-value pair.

The second part, the LT, is performed in a single transaction. In this part the thread first verifies the nodes that were found in the setup part are still valid (i.e., they are still live), their successive nodes are still live, and the pointers from their predecessors point to them. If one of the conditions does not hold, the transaction is aborted, and the whole remove operation is restarted. It then

Leap-List Update - Release and Update

```

input : size s, nodes pa, nodes na, nodes n, nodes new_node,
        booleans split
115 foreach  $j < s$  do
116   if  $split[j]$  then
117     if  $new\_node[j][1] \rightarrow level > new\_node[j][0] \rightarrow level$ 
        then
118       foreach  $i < new\_node[j][0] \rightarrow level$  do
119          $new\_node[j][0] \rightarrow next[i] := new\_node[j][1];$ 
120          $new\_node[j][1] \rightarrow next[i] :=$ 
          UNMARK( $n[j] \rightarrow next[i]$ );
121       end
122       foreach
         $new\_node[j][0] \rightarrow level \leq i < old\_node[j][1] \rightarrow level$ 
        do
123          $new\_node[j][1] \rightarrow next[i] :=$ 
          UNMARK( $n[j] \rightarrow next[i]$ );
124       end
125     else
126       foreach  $i < new\_node[j][1] \rightarrow level$  do
127          $new\_node[j][0] \rightarrow next[i] := new\_node[j][1];$ 
128          $new\_node[j][1] \rightarrow next[i] :=$ 
          UNMARK( $n[j] \rightarrow next[i]$ );
129       end
130       foreach
         $new\_node[j][1] \rightarrow level \leq i < old\_node[j][0] \rightarrow level$ 
        do
131          $new\_node[j][0] \rightarrow next[i] :=$ 
          UNMARK( $na[j][i]$ );
132       end
133     end
134   else
135     foreach  $i < new\_node[j][0] \rightarrow level$  do
136        $new\_node[j][0] \rightarrow next[i] :=$ 
        UNMARK( $n[j] \rightarrow next[i]$ );
137     end
138   end
139   foreach  $i < new\_node[j][0] \rightarrow level$  do
140      $pa[j][i] \rightarrow next[i] := new\_node[j][0];$ 
141   end
142   if  $split[j] \wedge (new\_node[j][1] \rightarrow level >$ 
         $new\_node[j][0] \rightarrow level)$  then
143     foreach
         $new\_node[j][0] \rightarrow level \leq i < old\_node[j][1] \rightarrow level$ 
        do
144        $pa[j][i] \rightarrow next[i] := new\_node[j][1];$ 
145     end
146   end
147    $new\_node[j][0] \rightarrow live := \mathbf{true};$ 
148   if  $split[j]$  then  $new\_node[j][1] \rightarrow live := \mathbf{true};$ 
149 end

```

Figure 10: Leap-List Update - Release and Update.

continues to mark the next pointers of the nodes that are about to be removed, and the next pointers of their predecessors. The transaction concludes by setting the live bit of the nodes to false, and attempts to commit. In case the commit fails, the remove operation is retried from the beginning of the setup part.

However, if the transaction successfully commits, the third part releases and updates each Leap-List to include the new nodes. It first sets the next pointers of the new node to point to the unmarked removed nodes next pointers in lines 217-227. Following this we

Leap-List Remove - Setup

```

input : Leap-Lists ll, keys k, values v, size s, nodes pa,
        nodes na, nodes n, nodes old_node, booleans merge,
        booleans changed
150 foreach  $j < s$  do
151   int total;
152   retry_last:  $merge[j] := \mathbf{false};$ 
153    $(pa, na) \leftarrow \text{PredecessorSearch}(ll[j], k[j]);$ 
154    $old\_node[j][0] := na[j][0];$ 
155   if  $get\_index(old\_node[j][0] \rightarrow trie, k[j]) =$ 
         $NOT\_FOUND$  then
156      $changed[j] := \mathbf{false};$ 
157     continue;
158   end
159   repeat
160      $old\_node[j][1] := old\_node[j][0] \rightarrow next[0];$ 
161     if  $\neg$  then goto  $retry\_last;$ 
162     until  $\neg is\_marked(old\_node[j][1]);$ 
163     total :=  $old\_node[j][0] \rightarrow count;$ 
164     if  $old\_node[j][1]$  then
165       total +=  $old\_node[j][1] \rightarrow count;$ 
166       if total  $\leq node\_size$  then  $merge[j] := \mathbf{true};$ 
167     end
168     Set  $n[j]$  level, count, high and low;
169     if  $\neg old\_node[j][0] \rightarrow live$  then goto  $retry\_last;$ 
170     if  $merge[j] \wedge \neg old\_node[j][1] \rightarrow live$  then goto  $retry\_last;$ 
171      $changed[j] := \text{RemoveAndMerge}(old\_node[j], n[j], k[j],$ 
         $merge[j]);$ 
172 end

```

Figure 11: Leap-List Remove - Setup.

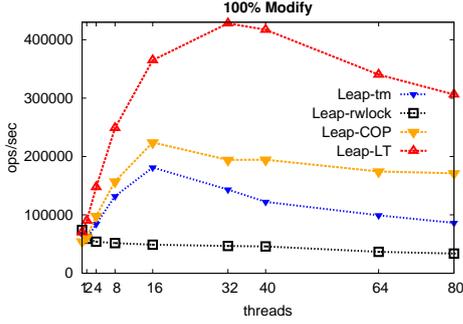
set the next pointers of the old nodes pointers to the new node (lines 229-230). It concludes, in line 232, by setting the new nodes live bit.

3. Evaluation

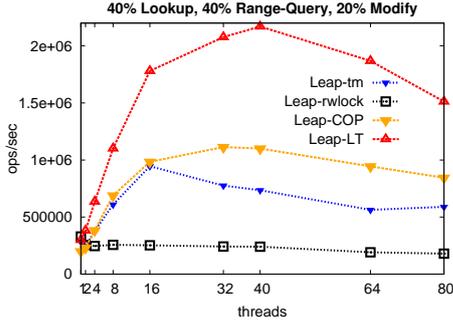
In this section we present the evaluation of our *Leap-List* implementation using COP and the LT technique and compare it to an STM-based *Leap-List*, an STM based *Leap-List* implementation that uses only COP, and a RW-Lock *Leap-List* implementation that uses a reader-writer lock. In Section 3.1 we compare to *Skip-list* implementations.

Experimental setup: We collected results on a machine powered by four Intel E7-4870. An Intel E7-4870 is a chip multithreading (CMT) processor, with 10 2.4 GHz cores each multiplexing 2 hardware threads, for a total of 20 hardware strands per chip. All implementations were compiled using GCC version 4.7 [1] which has built-in support for transactional memory. We used the linearizable memory allocation manager which was proposed in [7]. We compared the throughput (operations per second) of the following four algorithms:

1. **Leap-LT** - our proposed algorithm that uses COP and the LT technique as described in Section 2.
2. **Leap-tm** - a *Leap-List* implementation which wraps each operation within a transaction.
3. **Leap-COP** - an STM-based *Leap-List* implementation that uses COP (separating the search and update/remove operation).
4. **Leap-rwlock** - A Read-Write lock *Leap-List* implementation, in which the lookup and range-query operations acquire the

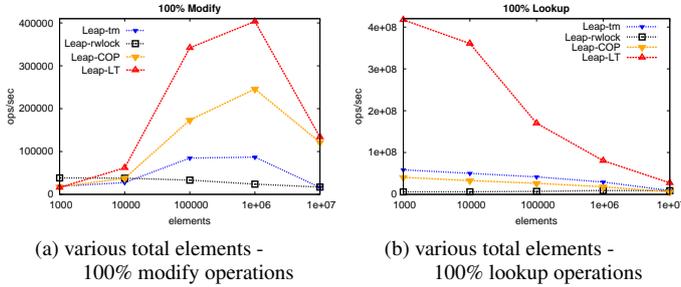


(a) various threads - 100% modify operations



(b) various threads - 40% lookup, 40% range-query, 20% modify operations

Figure 14: Leap-List size 100K. Workload: different amount of modifications (updates and removes), lookups and range queries. (a) 100% modify operations, (b) 40% lookup, 40% range-query and 20% modify operations.



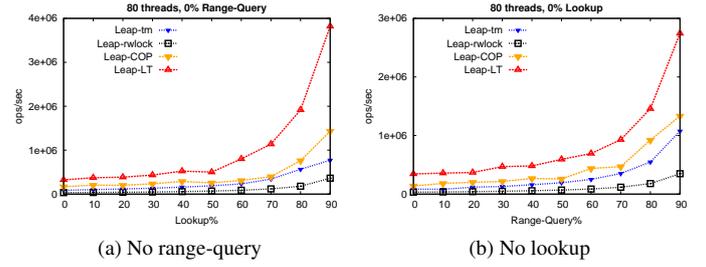
(a) various total elements - 100% modify operations

(b) various total elements - 100% lookup operations

Figure 15: Various total elements number. Workload: different amount of modifications (updates and removes) and lookups. (a) 100% modifications, (b) 100% lookups.

(compared to the lookup operations), and increased number of conflicts and retries.

Figure 15 shows the performance of the algorithms while varying the number of elements each *Leap-List* is initialized with, and setting the number of threads to 80. (The x-axis is log-scaled). We observe that when there are only update and remove operations (Figure 15-(a)), the highest throughput is achieved when a *Leap-List* is initialized with 1,000,000 elements. This is because there are less conflicts due to the high number of nodes. Note that when the number of elements is higher, the overhead stems from the long predecessors search operation. In figure 15-(b) we see that when there are only lookup operations, the highest throughput is achieved when the number of elements is 10,000. This is again due to the



(a) No range-query

(b) No lookup

Figure 16: Leap-List size 100K, 80 threads. Workload: different rates of modifications. (a) 0%-90% lookup and modify operations (no range-query), (b) 0%-90% range-query and modify operations (no lookup).

long predecessors search operations when the number of nodes is larger.

Figure 16-(a) and figure 16-(b) depict the throughput when using 80 threads, a *Leap-List* with 100,000 elements and varying the rate of lookup and range-query operations respectively between 0% to 90%. Both figures show that as the modifications rate is decreased, the throughput of all algorithms increases. In the case where no range-query operations occur (Figure 16-(a)) *Leap-LT* shows between 190% (0% lookup rate) to 260% (90% lookup rate) higher throughput compared with *Leap-COP*. The case where no lookup operations occur (Figure 16-(b)) exhibits similar results where *Leap-LT* shows between 240% (0% range-queries rate) to 200% (90% range-queries rate) higher throughput compared with *Leap-COP*. Note that in the case of 100% lookup and range-query operations rate (not shown here) the *Leap-LT* results are even better. *Leap-LT* is better by 650% and 320% compared to the second best *Leap-COP* in the 100% lookup and 100% range-query cases respectively.

3.1 Comparison to skip-lists

It is natural to compare our *Leap-LT* to the known *Skip-List* data-structure. We compare the throughput of various settings of a single *Leap-List* to: (1) *Skip-tm* - a skip-list implementation that uses the GCC-TM to synchronize operations; (2) *Skip-cas* - a skip-list implementation as described in [8]. These implementations store a single key-value pair in each node, and use mutable objects, thus having a lower modify operations overhead compared to our *Leap-LT*. Note that for this comparison we used a single *Leap-List* data-structure ($L = 1$), and that the range-query operations of the *Skip-cas* implementation do not return a consistent range-query (i.e., this operation is non-atomic and may return an inconsistent result).

Figures 17-(a), 17-(b), and 17-(c) show the throughput when using a data-structure with 1,000,000 elements, and varying the number of threads between 1 to 80. When there are only modify operations (Figure 17-(a)), we observe that both *Skip-cas* and *Skip-tm* are better than *Leap-LT*, and that *Skip-cas* is much better. This is due to the higher overhead of the update and remove operations in *Leap-LT*.

However, we see different results when there are more lookup and range-query operations, as can be seen in Figure 17-(b) where there are 40% lookups, 40% range-queries and 20% modifications. Here we see that *Leap-LT* is up to 2x and 38x better than *Skip-cas* and *Skip-tm* respectively. This is due to the overhead of the range-query operation that needs to iterate many nodes and to the large number of elements which reduces conflicts between concurrent modifying operations.

A workload which exhibits only lookup operations (Figure 17-(c)), shows that *Leap-LT* and *Skip-cas* are comparable and are much better than *Skip-tm*. This is because no contention occurs,

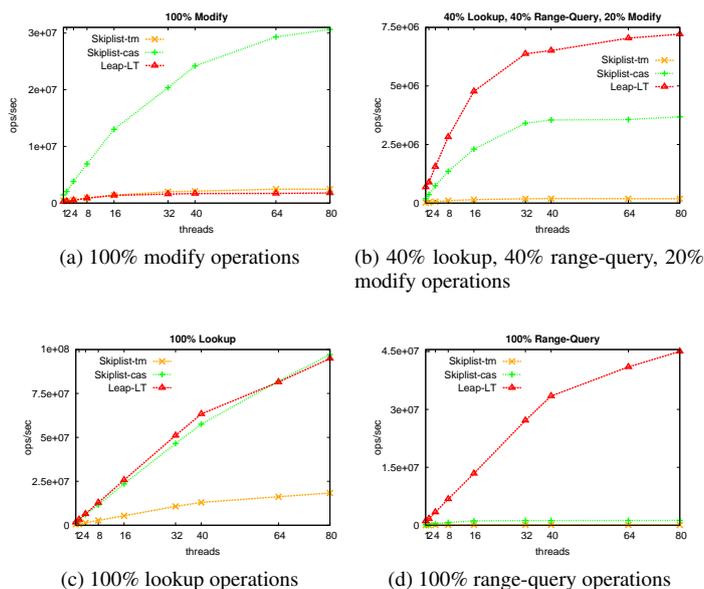


Figure 17: Leap-List comparison to Skip-Lists with 1M elements. Workload: different amount of modifications (updates and removes), lookups and range queries. (a) 100% modify operations, (b) 40% lookup, 40% range-query and 20% modify operations, (c) 100% lookup operations, (d) 100% range-query operations.

and the reduced overhead of the former algorithms produces better throughput.

Figure 17-(d) shows the main strength of our *Leap-LT* implementation on a workload of only range-query operations. It achieves better scalability and up to 35x better throughput on this workload compared to the *Skip-cas* implementation. Moreover, we note that this is achieved while ensuring a consistent operation result (which is not ensured in *Skip-cas*).

4. Summary

In this paper we presented a novel concurrent data-structure, *Leap-List*, that provides linearizable range queries. We implemented it using a technique called *Locking Transactions*, which reduces the executed transactions' lengths. We compared different *Leap-List* implementations, and also compared our technique to a *Skip-List* implementation.

We believe that the availability of hardware transactions will greatly enhance *Leap-List* performance because its implementation is based on short transactions. In the future we plan to test the *Leap-List* in an In-Memory Data-Base implementation, to replace the B-trees for indexes. We believe this can significantly improve the throughput of many Data-Base workloads.

5. Acknowledgements

This work was supported in part by NSF grant CCF-1217921 and by grants from the Oracle and Intel corporations.

References

[1] Gcc version 4.7.0, (<http://gcc.gnu.org/gcc-4.7/>), Apr. 2012. URL <http://gcc.gnu.org/gcc-4.7/>.
 [2] Y. Afek, H. Avni, and N. Shavit. Towards consistency oblivious programming. In *OPODIS*, pages 65–79, 2011.

[3] A. Braginsky and E. Petrank. A lock-free b+tree. In *SPAA*, pages 58–67, 2012.
 [4] N. G. Bronson, J. Casper, H. Chafi, and K. Olukotun. Transactional predication: high-performance concurrent sets and maps for stm. In *PODC*, pages 6–15, 2010.
 [5] T. Brown and H. Avni. Range queries in non-blocking k-ary search trees. In *OPODIS*, 2012.
 [6] P. Ferragina and R. Grossi. The string b-tree: a new data structure for string search in external memory and its applications. *J. ACM*, pages 236–280, 1999.
 [7] K. Fraser. *Practical lock freedom*. PhD thesis, Cambridge University Computer Laboratory, 2003. Also available as Technical Report UCAM-CL-TR-579.
 [8] K. Fraser. Practical lock-freedom. Ph. D. dissertation, UCAM-CL-TR-579, Computer Laboratory, University of Cambridge, 2004.
 [9] M. Herlihy, Y. Lev, V. Luchangco, and N. Shavit. A simple optimistic skiplist algorithm. In *Proceedings of the 14th international conference on Structural information and communication complexity*, pages 124–138, 2007.
 [10] A. Prokopec, N. G. Bronson, P. Bagwell, and M. Odersky. Concurrent tries with efficient non-blocking snapshots. In *Proceedings of the 17th ACM SIGPLAN symposium on Principles and Practice of Parallel Programming*, PPoPP '12, pages 151–160, 2012. ISBN 978-1-4503-1160-1.
 [11] W. Pugh. Skip lists: A probabilistic alternative to balanced trees. In *WADS*, pages 437–449, 1989.
 [12] O. Rodeh. B-trees, shadowing, and clones. *Trans. Storage*, 3(4):2:1–2:27, Feb. 2008. ISSN 1553-3077.