

Practical Non-blocking Unordered Lists

Kunlong Zhang¹, Yujiao Zhao¹, Yajun Yang¹, Yujie Liu², and Michael Spear²

¹ Tianjin University

{zhangkl, zyj0131, yangyajun}@tju.edu.cn

² Lehigh University

{yul510, spear}@cse.lehigh.edu

Abstract. This paper introduces new lock-free and wait-free unordered linked list algorithms. The composition of these algorithms according to the fast-path-slow-path methodology, a recently devised approach to creating fast wait-free data structures, is nontrivial, suggesting limitations to the applicability of the fast-path-slow-path methodology. The list algorithms introduced in this paper are shown to scale well across a variety of benchmarks, making them suitable for use both as standalone lists, and as the foundation for wait-free stacks and non-resizable hash tables.

1 Introduction

Linked lists are fundamental data structures that are widely used both on their own and as building blocks for other data structures. While a sequential linked list is easy to implement, concurrent linked lists that achieve both strong progress guarantees and good performance are challenging to design [3, 7–9, 16, 19, 22, 24]. Herlihy [10] demonstrated the existence of universal constructions for wait-free concurrent objects, yet it remains an open problem whether all such objects can be made practical: wait-free data structures implemented from universal constructions [4, 6, 11] tend to incur significant overhead, increased time and space complexity, and/or static bounds on the size of the data structure. Although many lock-free concurrent implementations [5, 12, 20, 21] have been proposed for sequential data structures, practical wait-free versions are relatively rare [14, 23].

We introduce the first practical implementation of an unordered linked list that supports *wait-free* insert, remove, and lookup operations. The implementation is linearizable [13] and uses only a single-word compare-and-swap (CAS) primitive. Furthermore, the implementation does not require marking the lower bits of pointers [8]. Our implementation is built from a novel lock-free unordered list algorithm, where each insert and remove operation first linearizes by appending an intermediate “request” node at the head of the list, followed by a lazy search phase that computes the return value of the operation (which depends on whether the key value is already in the set); lookup operations have no side-effects on the shared memory. The implementation achieves scalable wait-freedom by adapting a technique originally designed for wait-free queues [14], and to further improve performance, we applied a recently-devised fast-path-slow-path methodology [15] to construct adaptive variants of our algorithm.

In this paper, we introduce the first practical wait-free unordered linked list, which is immediately usable in applications as-is, and can be employed in the creation of

wait-free non-resizable hash tables and stacks.¹ We discuss our experience and findings in applying the fast-path-slow-path methodology, identifying both strengths and limitations of the approach. In Section 2, we present background and related work. In Section 3 we present a lock-free unordered list algorithm that serves as the basis for the wait-free algorithm discussed in Section 4. We evaluate performance in Section 5. Section 6 concludes with guidelines for using the fast-path-slow-path methodology.

2 Related Work

The first lock-free list to require only atomic compare-and-swap (CAS) operations was developed by Valois [24], who employed a technique in which auxiliary nodes encoded in-progress operations. Harris [8] implemented a lock-free ordered list by using a pointer marking technique, in which a node is logically deleted by marking the least significant bit of its next pointer; the node is then physically removed from the list in a separate phase. Michael [16] improved memory reclamation in the Harris algorithm using hazard pointers [17]. Heller et al. [9] designed a lock-based linked list with wait-free lookup operations. Their wait-free technique can also be incorporated into the Harris-Michael algorithm to improve performance. Kogan and Petrank [14] proposed a wait-free queue implementation and a more efficient variant based on the fast-path-slow-path methodology [15] which composes the slower wait-free algorithm with a faster lock-free implementation [18]. Timnat et al. [23] designed a wait-free ordered linked list based on the fast-path-slow-path methodology, using the Harris-Michael algorithm as its fast path.

Subsequent efforts have contributed to our general understanding of lock-free list implementations, but have neither improved progress guarantees nor delivered superior performance to that attainable by combining the Harris, Michael, and Heller techniques. Fomitchev and Ruppert [7] presented a lock-free list with worst-case linear amortized cost. Attiya and Hillel [1] presented a lock-free doubly-linked list that relies on a double-compare-and-swap (DCAS) operation. Sundell and Tsigas [22] presented a lock-free doubly-linked list using only CAS. Braginsky and Petrank [2] presented the first lock-free unrolled linked list.

Herlihy [10, 11] presented the first universal construction to convert sequential objects to wait-free concurrent implementations. Fatourou and Kallimanis [6] provided a universal construction that can be used to implement highly efficient stacks and queues.

3 A Lock-Free Unordered List

We now present a lock-free unordered list algorithm, which serves as the basis for our wait-free implementation. The algorithm implements a set object, where the elements can be compared using an equality operator ($=$), even if they can not be totally ordered.

The list supports three operations: $\text{INSERT}(k)$ attempts to insert value k into the set and returns true (success) if k was not present in the set, and returns false otherwise. $\text{REMOVE}(k)$ returns true if it successfully removes value k from the set and returns false if k does not exist in the set. $\text{CONTAINS}(k)$ indicates whether k is contained by the set.

¹ Presentation of these algorithms is included in a companion technical report [25].

```

datatype NODE
  key   :  $\mathbb{N}$  // integer data field
  state :  $\mathbb{N}$  // INS, REM, DAT, or INV
  next  : NODE // pointer to the successor
  prev  : NODE // pointer to the predecessor
  tid   :  $\mathbb{N}$  // thread id of the creator

global variables
  head : NODE // initially nil

1 function INSERT( $k : \mathbb{N}$ ) :  $\mathbb{B}$ 
2    $h \leftarrow \text{new NODE}(k, \text{INS}, \text{nil}, \text{nil}, \text{threadid})$ 
3   ENLIST( $h$ )
4    $b \leftarrow \text{HELPINSERT}(h, k)$ 
5   if  $\neg \text{CAS}(\&h.\text{state}, \text{INS}, (b? \text{DAT} : \text{INV}))$  then
6     HELPREMOVE( $h, k$ )
7      $h.\text{state} \leftarrow \text{INV}$ 
8   return  $b$ 

9 function REMOVE( $k : \mathbb{N}$ ) :  $\mathbb{B}$ 
10   $h \leftarrow \text{new NODE}(k, \text{REM}, \text{nil}, \text{nil}, \text{threadid})$ 
11  ENLIST( $h$ )
12   $b \leftarrow \text{HELPREMOVE}(h, k)$ 
13   $h.\text{state} \leftarrow \text{INV}$ 
14  return  $b$ 

15 function CONTAINS( $k : \mathbb{N}$ ) :  $\mathbb{B}$ 
16   $\text{curr} \leftarrow \text{head}$ 
17  while  $\text{curr} \neq \text{nil}$  do
18    if  $\text{curr}.\text{key} = k$  then
19       $s \leftarrow \text{curr}.\text{state}$ 
20      if  $s \neq \text{INV}$  then
21        return  $(s = \text{INS}) \vee (s = \text{DAT})$ 
22       $\text{curr} \leftarrow \text{curr}.\text{next}$ 
23  return false

24 procedure ENLIST( $h : \text{NODE}$ )
25  while true do
26     $\text{old} \leftarrow \text{head}$ 
27     $h.\text{next} \leftarrow \text{old}$ 
28    if  $\text{CAS}(\&\text{head}, \text{old}, h)$  then
29      return

30 function HELPINSERT( $h : \text{NODE}, k : \mathbb{N}$ ) :  $\mathbb{B}$ 
31   $\text{pred} \leftarrow h$ 
32   $\text{curr} \leftarrow \text{pred}.\text{next}$ 
33  while  $\text{curr} \neq \text{nil}$  do
34     $s \leftarrow \text{curr}.\text{state}$ 
35    if  $s = \text{INV}$  then
36       $\text{succ} \leftarrow \text{curr}.\text{next}$ 
37       $\text{pred}.\text{next} \leftarrow \text{succ}$ 
38       $\text{curr} \leftarrow \text{succ}$ 
39    else if  $\text{curr}.\text{key} \neq k$  then
40       $\text{pred} \leftarrow \text{curr}$ 
41       $\text{curr} \leftarrow \text{curr}.\text{next}$ 
42    else if  $s = \text{REM}$  then
43      return true
44    else if  $(s = \text{INS}) \vee (s = \text{DAT})$  then
45      return false
46  return true

47 function HELPREMOVE( $h : \text{NODE}, k : \mathbb{N}$ ) :  $\mathbb{B}$ 
48   $\text{pred} \leftarrow h$ 
49   $\text{curr} \leftarrow \text{pred}.\text{next}$ 
50  while  $\text{curr} \neq \text{nil}$  do
51     $s \leftarrow \text{curr}.\text{state}$ 
52    if  $s = \text{INV}$  then
53       $\text{succ} \leftarrow \text{curr}.\text{next}$ 
54       $\text{pred}.\text{next} \leftarrow \text{succ}$ 
55       $\text{curr} \leftarrow \text{succ}$ 
56    else if  $\text{curr}.\text{key} \neq k$  then
57       $\text{pred} \leftarrow \text{curr}$ 
58       $\text{curr} \leftarrow \text{curr}.\text{next}$ 
59    else if  $s = \text{REM}$  then
60      return false
61    else if  $s = \text{INS}$  then
62      if  $\text{CAS}(\&\text{curr}.\text{state}, \text{INS}, \text{REM})$  then
63        return true
64    else if  $s = \text{DAT}$  then
65       $\text{curr}.\text{state} \leftarrow \text{INV}$ 
66      return true
67  return false

```

Fig. 1. A Lock-free Unordered List

3.1 Overview

Figure 1 presents the basic algorithm. The list is comprised of `NODE` objects, where each `NODE` stores a *key* value, a *next* pointer to the successor node, and a *state* field for coordinating concurrent operations. The *prev* and *tid* fields are reserved for the wait-free algorithm (Section 4). We maintain a global pointer *head* that points to the first element of the list. Elements are always inserted at the head position.

The key insight of the algorithm is to maintain a refinement mapping function that maps a linked list object (starting from node *h*) to an abstract set object $\text{AbsSet}(h)$:

$$\text{AbsSet}(h) \equiv \begin{cases} \emptyset & \text{if } h = \text{nil} \\ \text{AbsSet}(h.\text{next}) & \text{if } h.\text{state} = \text{INV} \\ \text{AbsSet}(h.\text{next}) \cup \{h.\text{key}\} & \text{if } h.\text{state} = \text{INS} \vee h.\text{state} = \text{DAT} \\ \text{AbsSet}(h.\text{next}) \setminus \{h.\text{key}\} & \text{if } h.\text{state} = \text{REM} \end{cases}$$

To maintain this property, an INSERT or REMOVE operation first places a node with an intermediate state (*INS* or *REM*) at the head of the list. Then it searches the list for the value being inserted or removed, removing logically deleted nodes along the way. Finally, it sets the intermediate node to a final state (*DAT* or *INV*).

In more detail, an INSERT operation allocates an *INS* node (h) and links it to the head of the list by invoking ENLIST (lines 2 - 3). It then invokes HELPIINSERT (line 4) to determine whether the insertion is effective, that is, to check whether the key is already present in the set. The return value of HELPIINSERT dictates the return value of the INSERT operation, as well as the final state of h (line 5): if the key was absent from the set, $h.\text{state}$ is set to *DAT*, and the insertion becomes effective; otherwise, $h.\text{state}$ is set to *INV*, indicating that the insertion failed due to the key already being present in the set, and h becomes a garbage node that will be physically removed by some subsequent operation. The update of $h.\text{state}$ must use a CAS instruction (line 5), since a concurrent REMOVE that deletes the same key may attempt to change $h.\text{state}$ concurrently. If the CAS fails, it means the key was deleted concurrently and the thread will invoke HELPREMOVE (lines 6 - 7) to help the deleting thread to clean up the list.

Similarly, a REMOVE operation starts by inserting a *REM* node at the head position (lines 10 - 11). The real work of removal is delegated to the HELPREMOVE operation (line 12), which traverses the list to delete the specified key and returns a boolean value indicating whether the key was found (and deleted). Then node h is set to the *INV* state (line 13), allowing some subsequent operation to remove it from the list.

The CONTAINS operation has no side effect on shared memory (it is read-only). The operation traverses the list to find the specified key and skips any *INV* nodes (lines 18 - 20). If a non-*INV* node with the specified key is encountered, the operation returns true (found) if the node is in state *DAT* or *INS* (line 21). Otherwise, the node is in *REM* state, which represents a REMOVE operation that can be thought of as having already deleted the key from the suffix of the list, and hence, the CONTAINS operation immediately returns false.

3.2 ENLIST Operation

Both INSERT and REMOVE use the ENLIST operation to insert a node at the head position. In the lock-free algorithm, ENLIST repeatedly performs a CAS operation (line 28), attempting to change *head* to point to h , until the CAS succeeds. However, this approach fails to provide *wait-freedom*: since the CAS operation at line 28 of a specific thread may fail repeatedly, for an unbounded number of times (due to contention), the thread may starve in the ENLIST operation and make no progress. In Section 4, we introduce a wait-free ENLIST implementation, and show the algorithm can be made wait-free without any change to the other parts.

3.3 Coordination Protocol

The core protocol of coordinating concurrency is encapsulated by the `HELPINSERT` and `HELPREMOVE` operations. The two operations share a similar code structure: each takes a pointer parameter h , which points to the node inserted by the prior `ENLIST` operation. In both operations, the thread traverses the list starting from h , and reacts to the different types of nodes it encounters.

As a common obligation of both operations, logically deleted nodes are purged during the traversal (lines 35 - 38 and lines 52 - 55). That is, once an `INV` node is encountered (pointed to by $curr$), the node is physically removed from the list by setting the predecessor's next pointer to the successor of $curr$. Note that since new nodes cannot be added to the list at any point other than the head, the problems that plague node removal in sorted lists do not apply. In particular, it is not possible that removing one node can inadvertently lead to a new arrival disappearing from the list. While it is possible for a removed node to re-appear in the list on account of conflicting writes to the next pointer, such a node will necessarily already be marked `INV`, and thus there will be no impact on the correctness of the list.

During the traversal, the $curr$ node is skipped if $curr.key \neq h.key$ (lines 39 - 41 and 56 - 58). Otherwise, we say the $curr$ node is a "related node" with respect to the current operation. There are three possibilities if $curr$ is a related node: $curr$ is a `DAT` node, an `INS` node, or a `REM` node. In the latter two cases, the related node was created by some concurrent `INSERT` or `REMOVE` operation. We call such operations "related operations".

In `HELPINSERT`, if a related `REM` node is encountered, there is a concurrent `REMOVE` operation finalizing a removal of the same key. Hence, the `HELPINSERT` returns true (success) immediately (lines 42 - 43), since the concurrent `REMOVE` operation ensures that the key is absent in the set. Otherwise (lines 44 - 45), if the related node is an `INS` node, then the related `INSERT` operation inserted the same key earlier (or is determining that the key already exists in the list) and the `HELPINSERT` operation must return false. Finally, if the related node is a `DAT` node, `HELPINSERT` returns false since the key already exists in the set.

In `HELPREMOVE`, if a related `REM` node is found (lines 59 - 60), the operation returns false immediately since the key was already deleted by a concurrent `REMOVE` operation. If the related node is an `INS` node (lines 61 - 63), then the key was inserted by a concurrent `INSERT` operation. In this case, the thread attempts to change the node from `INS` to `REM` (line 62); a CAS instruction is needed to prevent data races on the `state` field (i.e., line 5). In the last case, the related node is a `DAT` node, meaning that the key is in the set, and the node is deleted by setting its `state` to `INV` (line 65).

3.4 Lock-Freedom

To show that the algorithm is lock-free, we show that *some* operation completes when any thread executes a bounded number of local steps. We first notice that the `ENLIST` operation is lock-free: a thread's CAS at line 28 may fail only due to another thread performing a CAS and completing its `ENLIST` operation. Since `ENLIST` is invoked exactly once in each `INSERT` and `REMOVE`, for n threads, at least one list operation will complete if some thread fails the CAS for n times in its `ENLIST` operation.

To show that every HELPINSERT and HELPREMOVE operation terminates, it is sufficient to show the list is acyclic. There are three places where the *next* pointer of a node is changed: executing line 27 cannot form a cycle, since the node *h* is newly allocated and is not reachable from any other node; when a thread executes line 37 or line 54, *pred* is clearly always a predecessor of *succ* in some total order *R*, which can be defined as the order in which nodes are inserted to the list (by the CAS at line 28).

Since the size of the list is bounded by *E*, the total number of completed ENLIST operations, every HELPINSERT and HELPREMOVE operation finishes in $O(E)$ steps. Note that in HELPREMOVE, a thread never executes the CAS at line 62 twice on the same node: if the CAS failed, the *curr* node is turned into a final state (*DAT* or *INV*) and will cause the loop to exit or skip the node in the next iteration. Thus, for *n* threads, either a thread completes its own list operation in $O(n + E)$ local steps, or some other thread completes a list operation during this period of time.

3.5 Linearizability

Due to space constraints, a complete proof of linearizability is provided in a companion technical report [25]. We define the linearization point for each operation: An INSERT(*k*) or REMOVE(*k*) operation linearizes at the successful CAS at line 28 in ENLIST. A CONTAINS(*k*) linearizes at line 16 if $k \notin \text{AbsSet}(\text{head})$ when *p* executes this line. In cases where $k \in \text{AbsSet}(\text{head})$ when *p* executes this line, the CONTAINS(*k*) linearizes at line 16 if the operation returns true. If the operation returns false, we show that there exists a concurrent REMOVE(*k*) that linearizes after *p* executes line 16 and before *p*'s CONTAINS(*k*) returns. We let *p*'s CONTAINS(*k*) linearize *immediately after* the linearization point of this REMOVE(*k*). Note that multiple CONTAINS(*k*) operations may be required to linearize after the same REMOVE(*k*) operation, and any two of these CONTAINS(*k*) operations can be ordered arbitrarily.

4 Achieving Wait-Freedom

The major challenge of the wait-free list algorithm lies in the implementation of a wait-free ENLIST operation. In this section, we present a wait-free ENLIST implementation adapted from a wait-free enqueue technique introduced by Kogan and Petrank [14]. We also introduce an adaptive wait-free algorithm which allows applications to trade off between average latency and worst-case latency of operations.

4.1 Wait-Free ENLIST Implementation

The enqueue technique introduced by Kogan and Petrank [14] provides a wait-free approach to append nodes at the tail of a list, but it is not immediately available as a solution to the ENLIST problem where nodes are appended at the head position. We employ *prev* fields to solve this problem. The additional code for implementing a wait-free ENLIST is presented in Figure 2.

```

datatype DESC
  phase      : ℕ           // integer phase number
  pending    : ℬ           // whether operation is pending
  node       : NODE       // pointer to the enqueueing node

global variables
  head       : NODE
  dummy      : NODE
  counter    : ℕ
  status     : DESC[THREADS]

initially
  head ← new NODE(-1, REM, nil, nil, -1)
  dummy ← new NODE(-, -, -, -, -)
  counter ← 0
  foreach d in status do
    d ← new DESC(-1, false, nil)

68 procedure ENLIST(h : NODE)
69   phase ← F&I(&counter)
70   status[threadid] ← new DESC(phase, true, h)
71   for tid ← 0 ... (THREADS - 1) do
72     HELPENLIST(tid, phase)
73   HELPFINISH()

74 function ISPENDING(tid : ℕ, phase : ℕ) : ℬ
75   d ← status[tid]
76   return d.pending ∧ (d.phase ≤ phase)

77 procedure HELPENLIST(tid : ℕ, phase : ℕ)
78   while ISPENDING(tid, phase) do
79     curr ← head
80     pred ← curr.prev
81     if curr = head then
82       if pred = nil then
83         if ISPENDING(tid, phase) then
84           n ← status[tid].node
85           if CAS(&curr.prev, nil, n) then
86             HELPFINISH()
87             return
88       else
89         HELPFINISH()

90 procedure HELPFINISH()
91   curr ← head
92   pred ← curr.prev
93   if (pred ≠ nil) ∧ (pred ≠ dummy) then
94     tid ← pred.tid
95     d ← status[tid]
96     if (curr = head) ∧ (pred = d.node) then
97       d' ← new DESC(d.phase, false, d.node)
98       CAS(&status[tid], d, d')
99       pred.next ← curr
100      CAS(&head, curr, pred)
101      curr.prev ← dummy

```

Fig. 2. A Wait-free ENLIST Implementation

The basic idea of the wait-free ENLIST algorithm is to let different ENLIST operations help each other to complete. The helping mechanism must ensure that every ENLIST operation reaches the response point in bounded number of steps (wait-freedom). This requires every thread to announce its intention by creating a descriptor entry in a *status* array before starting an operation. During its operation, the thread must visit each entry in the status array, helping other threads to make progress. To prevent starvation, each operation is assigned a *phase* number from a strictly increasing counter, and an operation only helps those with smaller phase numbers.

The wait-free ENLIST operation goes through six steps, as depicted in Figure 3:

- (a) The thread first announces its operation by creating a descriptor entry in its slot (indexed by its thread id) in the *status* array (line 70). The descriptor contains the *phase* number of the operation, a boolean *pending* field that indicates whether the operation is incomplete, and a pointer to the enlisting node. Once the descriptor is announced, the subsequent steps can be performed by the thread itself or by some helper thread.
- (b) The thread finds the node pointed to by *head*, and attempts to change its *prev* field to the enlisting node *h* by a CAS instruction (line 85).
- (c) The thread sets the *pending* flag of the operation descriptor to false by installing a new descriptor (line 98); this prevents concurrent helpers from retrying after the node is enlisted.

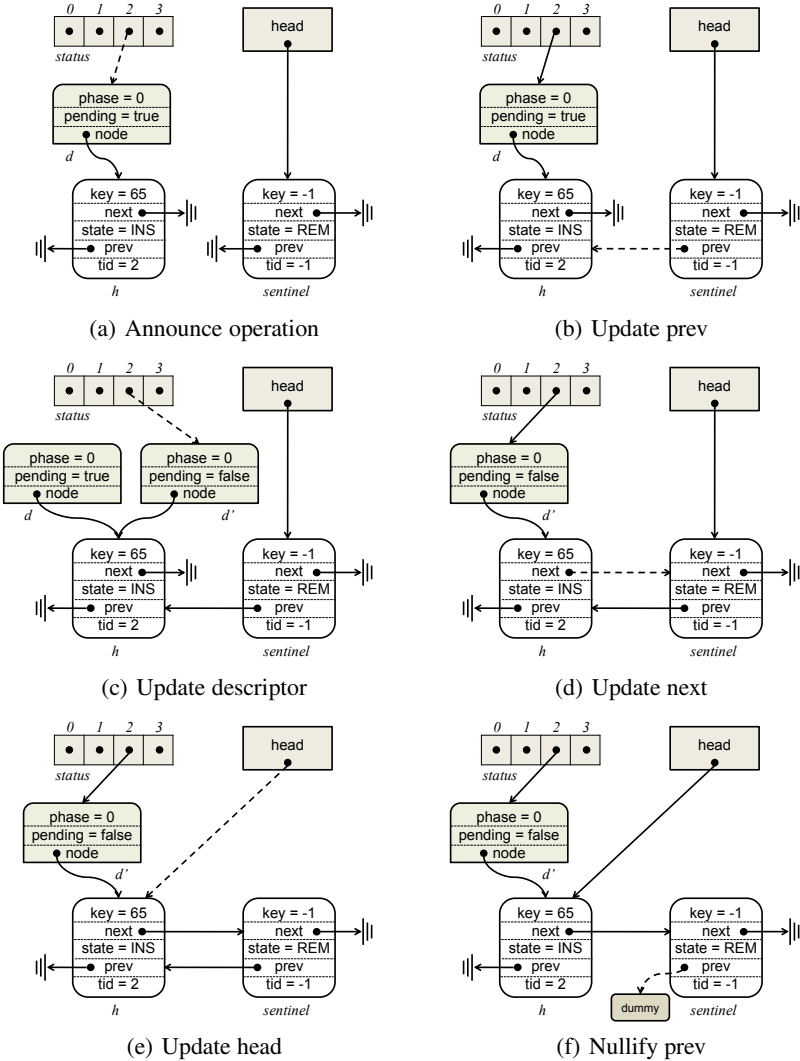


Fig. 3. Wait-free ENLIST Implementation Extended from the Kogan-Petrank Algorithm

(d) The thread sets $h.next$ to point to the original head node (line 99), which is the linearization point of the ENLIST operation. The ordering of this step is important with respect to steps (b) and (e). That is, the update of $h.next$ must be ordered after $head.prev$ is set to h , since the correct successor of h is “unknown” until then. On the other hand, $h.next$ must be updated before $head$ is changed to h , since otherwise a concurrent CONTAINS operation may start traversing from h and erroneously end by discovering $h.next$ is nil .

(e) The thread fixes $head$ by changing it to h using a CAS (line 100).

- (f) Finally, the thread clears the *prev* field of the original head by setting it to a *dummy* state (line 101). This is necessary for allowing the garbage collector to recycle deleted nodes. Since the *prev* pointers are installed by the wait-free ENLIST implementation, and the lock-free algorithm is unaware of their existence, keeping the *prev* pointers prevents the garbage collector from reclaiming a node even if the node is considered “unreachable” by the lock-free algorithm. It is worth noting that we must invalidate the *prev* pointer by setting it to a *dummy* state instead of **nil**, since the latter would admit ABA problems for the CAS instruction (line 85). Once the *prev* field of a node is set to *dummy*, it never changes.

4.2 An Adaptive Algorithm

Although the wait-free algorithm provides an upper bound on the steps required to complete an operation in the worst case, it imposes overhead in the common cases when contention is low. We employed the “fast-path-slow-path” methodology [15] to construct an adaptive algorithm that performs competitively in the common case while retaining the wait-free guarantee.

In the adaptive algorithm, a thread starts by executing a fast path version of the ENLIST operation, and falls back to the wait-free slow path if the fast path fails too many times (bounded by constant F). To prevent a thread from repeatedly taking the fast path while another thread starves, every thread checks the global status array after completing D operations, and performs helping if necessary. As shown in [15], for n threads, the adaptive algorithm ensures that every ENLIST operation completes in $O(F + D \cdot n^2)$ local steps. The F and D parameters can be adjusted to balance between the worst-case and common-case latency of operations.

It is worth noting that the fast path ENLIST of the adaptive algorithm is *not* equivalent to the lock-free ENLIST implementation in Figure 1. Instead, the fast path algorithm resembles the wait-free protocol, but excluding the announcing and helping steps.

5 Performance Evaluation

We evaluate performance of the lock-free and wait-free list algorithms via a set of microbenchmarks. These experiments allow us to vary the ratio of INSERT, REMOVE and CONTAINS operations, the range of key values, and the initial size of the list. We compare the following list-based set algorithms:

HarrisAMR: Implementation of the Harris-Michael algorithm [16] which also incorporates the wait-free CONTAINS technique introduced in [9]. The implementation uses Java `AtomicMarkableReference` objects to atomically mark deleted nodes.

HarrisRTTI: Optimized implementation of HarrisAMR in which Java run-time type information (RTTI) is used in place of `AtomicMarkableReference`. This is the best-known lock-free list implementation.

LazyList: Lock-based optimistic list implementation proposed by Heller et al [9].

LFList: The lock-free unordered list algorithm discussed in Section 3.

	Harris	LazyList	LFList	WFList	Adaptive
INSERT Cost	1 CAS	2 CAS	2 CAS	4 CAS + 1 F&I	3 CAS
REMOVE Cost	2 CAS	2 CAS	1 CAS	3 CAS + 1 F&I	2 CAS
Traverse Distance	$\frac{1}{2}k$		$(1 - \frac{\alpha}{2})k$		

Fig. 4. Update Cost and Average Traversal Distance (in uncontended cases)

WFList: The basic wait-free unordered list algorithm discussed in Section 4.

Adaptive: The adaptive wait-free unordered list algorithm discussed in Section 4.2.

FastPath: The fast-path portion of the Adaptive algorithm from Section 4.2.

In all implementations (except “HarrisAMR”), we use Java “FieldUpdaters” to perform CAS instructions on object fields. This approach provides better performance than simply using atomic fields (i.e. `AtomicInteger` and `AtomicReference`), which require expensive heap allocation cost and extra indirection overhead.

Experiments were conducted on an HP z600 machine with 6GB RAM and a 2.66GHz Intel Xeon X5650 processor with 6 cores (12 total threads) running Linux kernel 2.6.37 and OpenJDK 1.6.0. Each data point is the median of five 5-second trials. Variance was always below 5%.

5.1 Expected Overheads

Figure 4 enumerates the expected overheads of each of the algorithms. The cost of a successful list operation is affected by the update cost and the traversal cost. We measure the cost of an update operation (INSERT or REMOVE) by the number of atomic instructions required in the uncontended case. Compared to the Harris algorithm, LFList uses an extra CAS instruction in INSERT and one less in the REMOVE operation. The WFList requires 2 more CAS instructions to provide wait-freedom, though this cost is reduced in the Adaptive algorithm by leveraging the lock-free fast path.

The traversal cost is the average number of nodes that must be accessed. Suppose the list contains k elements uniformly selected from range $[0 \dots M)$ and let $k = \alpha M$ ($0 \leq \alpha \leq 1$). The average traversal distance for searching a random key value in an ordered list is: $D_o = \frac{1}{2}k$. In unordered lists, the average traversal distance is averaged among successful and unsuccessful search operations: $D_u = \alpha \cdot \frac{1}{2}k + (1 - \alpha)k = (1 - \frac{\alpha}{2})k$. This suggests that ordered lists have an increasing advantage over unordered lists when the set is sparse. For instance, when $\alpha = \frac{1}{2}$ (half of key space is in the set), the average traversal distance in an unordered list is 50% longer than its ordered permutation. Note too that in the ordered lists, an unsuccessful insert/remove does not perform a CAS, whereas every insert/remove in the unordered list performs a CAS.

5.2 x86 Performance

In Figures 5–7, we assess the performance of the lists for a variety of workloads. The “L” parameter indicates the percentage of operations that are lookups, with the

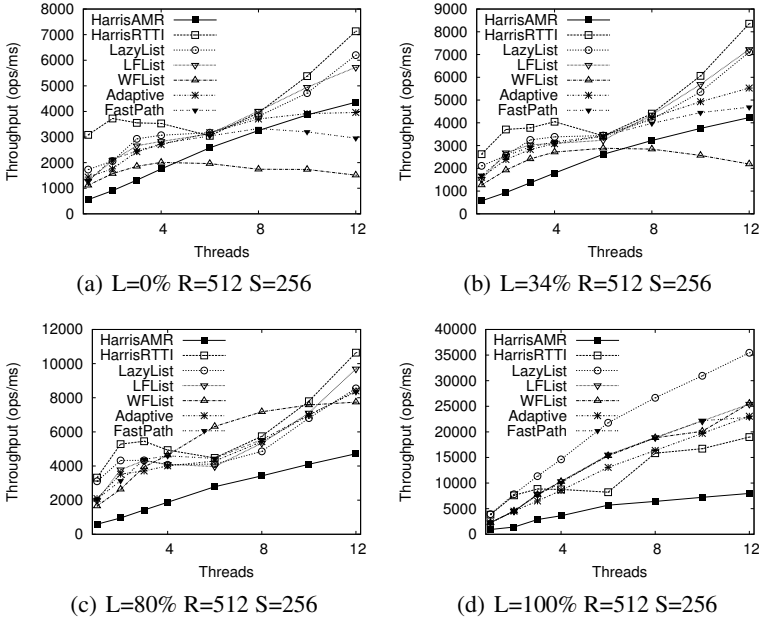


Fig. 5. Microbenchmark - Short Lists (L: Lookup Ratio, R: Key Range, S: List Size)

remainder evenly split between inserts and removals. “R” indicates the key range, and “S” indicates the average size of the list. In every case, the list is pre-populated with a random selection of S unique elements in the range [0, R). These elements are chosen at random, without replacement. Thus in the unordered lists, they will not be ordered.

The x86 processor features an aggressive pipeline, a deep cache hierarchy, and low-latency CAS operations. On this platform, the cost of write-write sharing is high, and thus both the wait-free enlistment mechanism and conflicting CAS operations on the head of the list are potential scalability bottlenecks. Nonetheless, our lock-free and wait-free algorithms scale well in all but a few cases. Indeed, the difference in performance appears to be much more a consequence of the increased traversal distance in the unordered algorithm than a consequence of increased cache misses due to frequent updates to the head of the list.

The most immediate and consistent finding is that the Harris list without RTTI optimizations has substantially higher latency and worse scalability than all other algorithms. We include this result as a reminder that concurrent data structures must be implemented using state-of-the-art techniques. Merely showing improved performance relative to the canonical Harris list presented in [12] does not give any indication of real-world performance. In particular, we caution that a direct comparison between our list and the wait-free ordered list [23] is not possible until that list is redesigned to use these modern optimizations.

We also see that long-running and read-only operations significantly reduce the cost of wait-free enlistment. When lists are small and updates are frequent, the enlistment

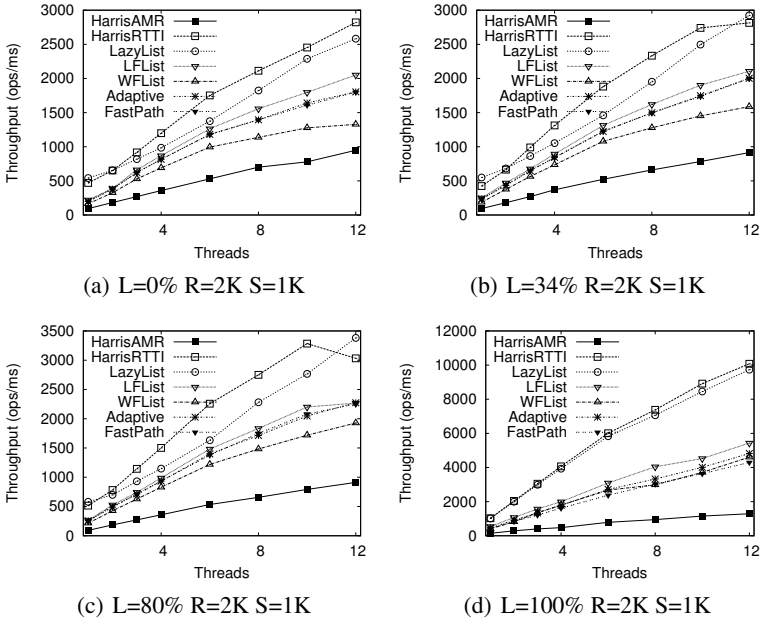


Fig. 6. Microbenchmark - Medium Lists (L: Lookup Ratio, R: Key Range, S: List Size)

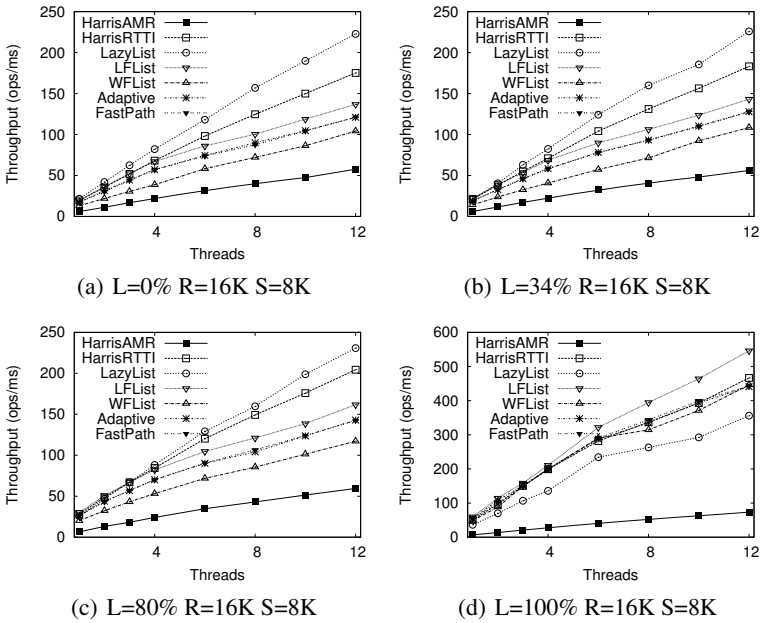


Fig. 7. Microbenchmark - Long Lists (L: Lookup Ratio, R: Key Range, S: List Size)

table and counter themselves become a bottleneck. Otherwise, the adaptive algorithm and its FastPath component are nearly identical.

The FastPath lock-free list is always a constant factor slower than the lock-free unordered list, but the Adaptive algorithm remains close to FastPath. This finding confirms Kogan and Petrank's claim [15] that the fast-path-slow-path technique can provide worst-case wait-freedom with lock-free performance. Furthermore, since the average operation in our list accesses many locations, contention on the head node of the list, while significant, does not dominate. Thus we observed that even for small thresholds, the adaptive algorithm rarely fell back to wait-free mode. However, it is important to observe that the lock-free FastPath algorithm itself is slower than our best lock-free unordered list. We shall return to this point in Section 6.

6 Discussion and Future Work

In their paper introducing the fast-path-slow-path methodology, Kogan and Petrank state that "... each operation is built from a fast path and a slow path, where the former is a version of a lock-free implementation of that operation, and the latter is a version of a wait-free implementation. Both implementations are customized to cooperate with each other [15, Sec. 3]."

Given a lock-free algorithm L , the question then is how to apply the methodology to produce a wait-free algorithm that does not sacrifice performance. We will consider L as consisting of three phases: a prefix (instructions that occur before the linearization point), a CAS operation (the linearization point), and a suffix (clean-up operations that follow the linearization point). Considering the three existing fast-path-slow-path algorithms (this work, ordered lists [23], and queues [15]), we see a pattern emerge.

First, a correct wait-free algorithm W must be constructed. This entails adding an announcement operation and operation descriptors to L . However, this step introduces the possibility of helping in the prefix, and thus makes it possible for helping operations to race (particularly if there are stores to memory that would not be shared in L). To correct these races, extra fields must be added to nodes of the data structure, stores must be upgraded to CAS instructions, and these CAS instructions must be sequenced by performing intermediate updates (via CAS) to a descriptor after each prefix step. It appears that changes to the suffix of the operation are not required, since the suffix is either clean-up operations that already support helping (e.g., the second CAS in the M&S queue [18]), or else operations that do not affect data structure invariants (e.g., the list traversal in HELPINSERT).

The second step is to perform a reduction that yields a lock-free algorithm L' that remains compatible with W . The first step of the reduction is to elide the announce operation and descriptor updates in L' . Then W must be analyzed, step-by-step, and simplified in an ad-hoc manner. In the ideal case, the result is the original lock-free algorithm L . Currently, it appears that the ideal case only occurs when the prefix is empty and the linearization point is the first CAS. Otherwise (as is the case in our list and the ordered list [23]), L' will need additional CAS instructions (relative to L) to keep its prefix compatible with the prefix of W .

Nonetheless, the ability to create low-latency wait-free data structures is valuable, particularly data structures as fundamental as linked lists. To emphasize the

significance of our wait-free unordered list, note that our list can be extended to support a REMOVEHEAD operation. Such an operation would resemble our REMOVE operation, but using a wildcard as its key value, and would immediately yield a wait-free stack. In contrast to stacks, constructing wait-free resizable hash tables based on our lists will be nontrivial. One challenge is that the shared descriptor array may become a bottleneck; were it not for resizing, each bucket could have its own descriptor array. However, the unordered nature may simplify other aspects of the design, for example, easing the implementation of list merging/splitting since the resulting lists need not be sorted.

Acknowledgements. We would like to thank Tim Harris, Alex Kogan, Victor Luchangco and our anonymous reviewers for their helpful suggestions during the preparation of our final manuscript.

References

1. Attiya, H., Hillel, E.: Built-In Coloring for Highly-Concurrent Doubly-Linked Lists. In: Dolev, S. (ed.) DISC 2006. LNCS, vol. 4167, pp. 31–45. Springer, Heidelberg (2006)
2. Braginsky, A., Petrank, E.: Locality-Conscious Lock-Free Linked Lists. In: Proceedings of the 12th International Conference on Distributed Computing and Networking, Bangalore, India (January 2011)
3. Colvin, R., Groves, L., Luchangco, V., Moir, M.: Formal verification of a lazy concurrent list-based set algorithm. In: Ball, T., Jones, R.B. (eds.) CAV 2006. LNCS, vol. 4144, pp. 475–488. Springer, Heidelberg (2006)
4. Ellen, F., Fatourou, P., Kosmas, E., Milani, A., Travers, C.: Universal Constructions that Ensure Disjoint-Access Parallelism and Wait-Freedom. In: Proceedings of the 2012 ACM Symposium on Principles of Distributed Computing (July 2012)
5. Ellen, F., Fatourou, P., Ruppert, E., van Breugel, F.: Non-blocking Binary Search Trees. In: Proceedings of the 29th ACM Symposium on Principles of Distributed Computing, Zurich, Switzerland (July 2010)
6. Fatourou, P., Kallimanis, N.D.: A Highly-Efficient Wait-Free Universal Construction. In: Proceedings of the 23rd ACM Symposium on Parallelism in Algorithms and Architectures, San Jose, CA (June 2011)
7. Fomitchev, M., Ruppert, E.: Lock-Free Linked Lists and Skip Lists. In: Proceedings of the 23rd ACM Symposium on Principles of Distributed Computing, St. John's, Newfoundland, Canada (July 2004)
8. Harris, T.L.: A Pragmatic Implementation of Non-Blocking Linked Lists. In: Welch, J.L. (ed.) DISC 2001. LNCS, vol. 2180, pp. 300–314. Springer, Heidelberg (2001)
9. Heller, S., Herlihy, M.P., Luchangco, V., Moir, M., Scherer III, W.N., Shavit, N.N.: A Lazy Concurrent List-Based Set Algorithm. In: Anderson, J.H., Prencipe, G., Wattenhofer, R. (eds.) OPODIS 2005. LNCS, vol. 3974, pp. 3–16. Springer, Heidelberg (2006)
10. Herlihy, M.: Wait-Free Synchronization. *ACM Transactions on Programming Languages and Systems* 13(1), 124–149 (1991)
11. Herlihy, M.: A Methodology for Implementing Highly Concurrent Data Objects. *ACM Transactions on Programming Languages and Systems* 15(5), 745–770 (1993)
12. Herlihy, M., Shavit, N.: *The Art of Multiprocessor Programming*. Morgan Kaufmann (2008)
13. Herlihy, M.P., Wing, J.M.: Linearizability: A Correctness Condition for Concurrent Objects. *ACM Transactions on Programming Languages and Systems* 12(3), 463–492 (1990)

14. Kogan, A., Petrank, E.: Wait-Free Queues with Multiple Enqueuers and Dequeuers. In: Proceedings of the 16th ACM Symposium on Principles and Practice of Parallel Programming, San Antonio, TX (February 2011)
15. Kogan, A., Petrank, E.: A Methodology for Creating Fast Wait-Free Data Structures. In: Proceedings of the 16th ACM Symposium on Principles and Practice of Parallel Programming, New Orleans, LA (February 2012)
16. Michael, M.: High Performance Dynamic Lock-Free Hash Tables and List-Based Sets. In: Proceedings of the 14th ACM Symposium on Parallel Algorithms and Architectures, Winnipeg, Manitoba, Canada (August 2002)
17. Michael, M.: Hazard Pointers: Safe Memory Reclamation for Lock-Free Objects. *IEEE Transactions on Parallel and Distributed Systems* 15(6), 491–504 (2004)
18. Michael, M.M., Scott, M.L.: Simple, Fast, and Practical Non-Blocking and Blocking Concurrent Queue Algorithms. In: Proceedings of the 15th ACM Symposium on Principles of Distributed Computing (May 1996)
19. O’Hearn, P.W., Rinetzky, N., Vechev, M.T., Yahav, E., Yorsh, G.: Verifying Linearizability with Hindsight. In: Proceedings of the 29th ACM Symposium on Principles of Distributed Computing, Zurich, Switzerland (July 2010)
20. Prokopec, A., Bronson, N., Bagwell, P., Odersky, M.: Concurrent Tries with Efficient Non-Blocking Snapshots. In: Proceedings of the 17th ACM Symposium on Principles and Practice of Parallel Programming (February 2012)
21. Sundell, H., Tsigas, P.: Fast and Lock-Free Concurrent Priority Queues for Multi-Thread Systems. *Journal of Parallel and Distributed Computing* 65, 609–627 (2005)
22. Sundell, H., Tsigas, P.: Lock-Free Deques and Doubly Linked Lists. *Journal of Parallel and Distributed Computing* 68(7) (July 2008)
23. Timnat, S., Braginsky, A., Kogan, A., Petrank, E.: Wait-Free Linked-Lists. In: Baldoni, R., Flocchini, P., Binoy, R. (eds.) *OPODIS 2012*. LNCS, vol. 7702, pp. 330–344. Springer, Heidelberg (2012)
24. Valois, J.: Lock-free linked lists using compare-and-swap. In: Proceedings of the Fourteenth Annual ACM Symposium on Principles of Distributed Computing, Ottawa, Ontario, Canada (August 1995)
25. Zhang, K., Zhao, Y., Yang, Y., Liu, Y., Spear, M.: Practical Non-blocking Unordered Lists. Technical Report LU-CSE-13-003, Lehigh University (2013)