

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/337020321>

# A C++ Implementation of a Lock-Free Priority Queue Based on Multi-Dimensional Linked List

Experiment Findings · November 2019

CITATIONS

0

READS

2,272

2 authors, including:



Steven Carroll

University of Central Florida

7 PUBLICATIONS 149 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Threadsafe and Lock Free Priority Queues Following Methods of Zhang and Dechev (2016) [View project](#)

# A C++ Implementation of a Lock-Free Priority Queue Based on Multi-Dimensional Linked List

Alexander Goponenko  
Department of Computer Science  
University of Central Florida  
Orlando, FL 32826  
agoponenko@Knights.ucf.edu

Steven Carroll  
Department of Computer Science  
University of Central Florida  
Orlando, FL 32826  
sdc@Knights.ucf.edu

**Abstract**—This paper is the second in a series of publications describing a reimplement of the priority queue based on multi-dimensional linked lists (Zhang and Dechev [1]). This variant of priority queue guarantees  $O(\log N)$  worst-case time complexity where  $N$  is the size of the key universe. It also provides improved performance over the state of the art approaches under high concurrency because each insertion modifies at most two consecutive nodes, allowing concurrent insertions to be executed with minimal interference. The current publication is dedicated to a lock-free priority queue. Alternative implementations are briefly discussed.

**Index Terms**—priority queue, multi-dimensional list

## I. INTRODUCTION

A priority queue is a fundamental data structure that comprises a set of key-value pairs where keys indicate priorities (by convention, a smaller key indicates higher priority). A typical priority queue implements only two operations: **Insert**, which adds an item with its associated priority to the queue, and **DeleteMin**, which removes the highest priority item from the queue. This data structure is employed abundantly in algorithms everywhere from high-level applications to low-level system kernels. Its efficient implementation in multithreaded environments is critical for modern and future multi-core systems. [1] introduced a lock-free priority queue based on multi-dimensional linked lists, with worst-case time  $O(\log N)$  for a key universe of size  $N$ . This paper is the second in a series of publications about a reimplement of the priority queue from [1]. The current publication describes a lock-free priority queue. Alternative approaches are also briefly described.

## II. EXISTING APPROACHES

Both sequential and concurrent priority queues are generally built from array-based heaps. A heap can be viewed as a binary tree. A min-heap maintains the min-heap property, which states that every node is smaller than both of its children. Array-based heaps display heavy memory contention when a newly inserted key ascends to its target location and after the top item is removed [1]. Furthermore, the heap invariant is not maintained most of the time during the execution of an **Insert** or **DeleteMin** operation, impeding concurrent access to the structure. Therefore, efficient lock-free heap-based implementations of a priority queue are lacking. The heap

based implementation from Intel Threading Building Blocks, an established industry standard concurrent library, employs a dedicated aggregator thread to perform all operations and therefore is not lock-free [1].

The structure of the Skiplist [2] enables lock-free priority queue implementations that outperform heap-based implementations on modern multi-core systems. The first lock-free priority queues based on skiplists were presented by Sundell and Tsigas [3] and Herlihy and Shavit [4]. A more recent and enhanced implementation of skiplist-based priority queue was published by Linden and Jonsson [5]. Although skiplists proved to be an efficient base for concurrent priority queue implementations, they have several potential drawbacks. For instance, although average time complexity of inserting an item in a skiplist is  $O(\log N)$ , the worst-case time is  $O(N)$ , where  $N$  is the number of items in the structure. Also, inserting an item into a skiplist involves updating several distant nodes, which may cause interference among concurrent operations and reduce throughput in a lock-free scenario.

The major bottleneck of concurrent priority queues is due to the inherent sequential semantics of the **DeleteMin** operation [6]. Two approaches have been employed to alleviate this bottleneck: semantics relaxation and consistency relaxation [1]. The first approaches were explored by Alistarh et al. [7], who introduced the Spraylists structure, Wimmer et al. [8], who developed  $k$ -log-structured merge-trees, and Rihani et al. [9], who pioneered MultiQueue. The second set of approaches was demonstrated by Herlihy and Shavit [4] and Zhang and Dechev [1] through implementation of quiescently consistent rather than linearizable priority queues.

## III. MULTI-DIMENSIONAL LISTS

As described in [1], a multi-dimensional (linked) list organizes data in multiple dimensions in a way that facilitates search and insertion. Upon insertion, a scalar key is recalculated into an array of indexes that may be considered multi-dimensional coordinates. This array is used, starting at the lowest dimension, to find pivot points into higher dimensions during the search for the correct insertion point. In contrast to a one-dimensional linked list, many intermediate values can be, and often are, effectively skipped during an insertion. The same is true of the search operation.

Algorithm 1: MDList structures

```

1  class MDList {
2      const int D;
3      const int N;
4      Node* head;
5  }
6
7  struct Node {
8      int key, k[D];
9      void* val;
10     Node* child[D];
11 }

```

Algorithm 2: Mapping from integer to vector

```

1  vector<int> keyToCoord(int key) {
2      int basis = ceil(pow(N, 1.0 / D));
3      int quotient = key;
4      vector<int> k;
5      k.resize(D);
6      for (int i = D - 1; i ≥ 0; i--) {
7          k[i] = quotient % basis;
8          quotient = quotient / basis;
9      }
10     return k;
11 }

```

More formally, A  $D$ -dimensional list is a rooted tree in which each node is implicitly assigned a dimension from 0 to  $D - 1$ . The root node's dimension is 0. A node of dimension  $d$  has no more than  $D - d$  children, and each of its children is assigned a unique dimension in the range from  $d$  to  $D - 1$ . The order among nodes is lexicographically based on keys. A dimension  $d$  node should share a coordinate prefix of length exactly  $d$  with its parent [1].

Each insertion or deletion operation on an MDList requires updating at most two consecutive nodes in the data structure, which makes it suitable for concurrent accesses. Furthermore, the worst-case time of the operations is  $O(\log N)$ , where  $N$  is the size of key universe.

Algorithms 1, 2, and 3 are reproduced from [1]. The priority queue is represented by a class called MDList which has a head field, a key universe size field,  $N$ , and a dimension field,  $D$ .

As mentioned in [1], the insert algorithm requires only a simple modification to the search algorithm, where a *prev* pointer is maintained for linking when the proper position is found.

In this work we will consider only integer values and keys.

Algorithm 3: Search for a Node with Coordinates

```

1  Node* searchNode(vector<int> k) {
2      Node *cur = head;
3      int d = 0;
4      while (d < D) {
5          while (cur ≠ NULL && k[d] > cur->k[d])
6              cur = cur->child[d];
7          if (cur == NULL || k[d] < cur->k[d])
8              return NULL;
9          d++;
10     }
11     return cur;
12 }

```

Algorithm 4: Data Structures for Lock-Free Priority Queue

```

1  struct AdoptDesc{
2      Node* curr;
3      int dp, dc;
4  };
5
6  struct Node{
7      atomic<uintptr_t> val;
8      int key;
9      int k[D];
10     atomic<uintptr_t> child[D];
11     atomic<AdoptDesc*> adesc;
12 };
13
14 struct HeadNode: Node{
15     int ver;
16 };
17
18 struct Stack{
19     Node* node[D];
20     HeadNode* head;
21 };

```

Algorithm 5: Marking Scripts

```

1  #define SetMark(p,m) ((p)|(m))
2  #define ClearMark(p,m) ((p)&~(uintptr_t)(m))
3  #define IsMarked(p,m) ((p)&(uintptr_t)(m))
4  #define F_ADP 0x1U
5  #define F_PRG 0x2U
6  #define F_DEL 0x1U
7  #define F_ALL 0x3U
8  #define Clear(p) ClearMark(p, F_ALL)

```

As in [1], only unique keys will be considered. Converting the key into an appropriate  $D$ -dimensional vector is accomplished by Algorithm 2.

#### IV. LOCK-FREE IMPLEMENTATION

Algorithm 4 shows the structure of the priority queue's node and other data structures. In addition to regular fields *key*, *k*, *child*, and *val*, the node contains field *adesc*, which holds a reference to a descriptor [4] if the node is just inserted and may not have adopted children yet (details are described later). The descriptor (*AdoptDesc*) contains a reference to the node from which the children must be adopted and a range of dimensions for which the adoption is pending. *HeadNode* has an additional field – *ver*, which holds the current version of the head. The version of the head is incremented during the purge operation that is described later.

The field *val* of the node has two purposes. Normally, it contains a reference to the node's value. However, if the node has been deleted, *val* is reused to hold references needed to maintain the queue after purge. The deleted node is marked using the so-called "bit stealing" technique – the last bit of *val* is set to 1 (in this case we will say that the flag *F\_DEL* is set). The same technique is used to mark invalid references in the *child* array. If the reference is invalid because it was adopted, it is marked with *F\_ADP*; if it is invalid because of a purge, it is marked with *F\_PRG*. Macros used to mark and unmark references are shown in Algorithm 5.

The structure of the priority queue itself is demonstrated in Algorithm 6.  $N$  contains the limit on the value of keys (the

### Algorithm 6: Priority Queue

```

1  class PriorityQueue {
2      int N;
3      int R
4      atomic_bool notPurging{true};
5      atomic<int> nMarkedNodes{0};
6      atomic<uintptr_t> head;
7      atomic<Stack*> stack
8      HeadNode firstHeadNode;
9      Stack firstStack;

11  PriorityQueue(int N, int R): N(N), R(R){
12      firstHeadNode.val = F_DEL;
13      firstHeadNode.adesc = NULL;
14      firstHeadNode.key = 0;
15      setCoords(&firstHeadNode, 0);
16      firstHeadNode.ver = 1;
17      for (int i=0; i<D; i++)
18          firstHeadNode.child[i].store(NIL);
19      head.store((uintptr_t) (&firstHeadNode));
20      firstStack.head = &firstHeadNode;
21      for (int i=0; i<D; i++)
22          firstStack.node[i] = &firstHeadNode;
23      stack.store(&firstStack);
24  }

26  void setCoords(Node* n, int key) {
27      int basis = ceil(pow(N, 1.0/D));
28      int quotient = key;
29      int* k = n->k;
30      for (int i = D - 1; i ≥ 0; i--) {
31          k[i] = quotient % basis;
32          quotient = quotient / basis;
33      }
34  }
35  };

```

size of the key universe);  $R$  holds the number of `deleteMin` operations between purges; `notPurging` is a flag needed to make sure that only one `purge` operation is ongoing; `nMarkedNodes` is the count of deleted nodes after the last purge; `head` is the reference to the current head node; `stack` contains the pointer to the current deletion stack. The constructor initializes a first head and a first stack of the queue. The head of the queue is a sentinel node, which is marked as deleted and has `key = 0`. The values of the keys of nodes in the queue must be between 0 and  $N$ . `setCoords` perform mapping from `key` to array `k`.

In order to attain lock-free and efficient functioning of the priority queue, the operations `deleteMin` and `insert` must be well coordinated. The nodes are deleted from the queue logically by setting the flag `F_DEL`. The deletion stack is used to store the location of the last logically deleted node to make physical deletion of nodes more efficient. A new node could be inserted in-between logically deleted nodes in a location that is not accessible from the current stack. Therefore, the `insert` operation may need to rewind the stack to make the newly inserted node accessible.

Insertion of a node is performed in two steps. In the first step the node is spliced into the list using a `CompareAndSwap` (CAS) atomic synchronization primitive in a way similar to that used in a lock-free linked list [10]. In the second step, the node adopts some of the children of the node that occupied its place if the dimension of the replaced node has been changed as a result of the insertion. The need for the second

step is announced by descriptor object. Other threads help the adoption if they traverse a node with a descriptor that has been set.

When the number of logically deleted nodes agglomerates above a threshold determined by the variable  $R$ , a `purge` operation is performed to ensure efficient execution and to enable memory reclamation. The `purge` operation may need to update the deletion stack and also must ensure that all non-deleted nodes remain accessible from the stack. The details of the operations are discussed next.

#### A. Details of insert operation

Algorithm 7 presents the operation of inserting an item into the priority queue. At the beginning of the operation a new node and a new stack are created. Memory allocation for new elements is dedicated to an object that implements a `Handler` interface, which is shown in Algorithm 8. The new stack gets recalculated throughout the operation in case the deletion stack needs to be rewound.

`LocatePlace` traverses the MDList starting from the head and determines the target position for the insertion, i.e. the immediate parent `pred` for the new node, dimension `dp`, at which the new node will become the child of `pred`, the node `curr` that is currently occupying the new node's slot, and dimension `dc`, at which `curr` will become the child of the new node. Nodes `pred` and `curr` are the only two nodes that are updated by an insertion. During traversal, `finishInserting` (Algorithm 9) is called for each inspected node, to complete possible ongoing adoption of children.

The CAS operation on line 21 splices the new node into the list. It can fail if the desired location has been updated by a concurrent insertion or because the location was marked invalid by a purge or a child adoption process. In such cases the loop beginning at line 8 restarts. Otherwise, the child adoption is completed (`finishInserting` at line 22) and the deletion stack is rewound if needed (`RewindStack` at line 23).

#### B. Details of deleteMin operation

Algorithm 10 demonstrates extraction of the item with highest priority, i.e. deletion of the node with the smallest key. The operation searches for a node that has not been logically deleted, starting from the last entry of the deletion stack. A copy of the stack is maintained following the search in order to update the queue's stack at the end of the operation. When a non-logically-deleted node is found, an attempt to logically delete it is performed by changing `val` to `F_DEL` with CAS atomic operation (line 27). In the case of success, the queue's stack is updated unless it has already been updated by another thread (line 30).

If the count of deleted nodes `nMarkedNodes` surpasses threshold  $R$ , a `purge` operation is attempted (line 33). Importantly, during the traversal of deleted nodes, `val` is inspected for the presence of a reference to a newer version of the queue's head (line 17). If a reference is found, the search continues from the newer head (lines 21-23).

Algorithm 7: Inserting a Node into MDList

```

1  bool insert(int key, uintptr_t val, Handler* h){
2      Stack *s = h->newStack();
3      Node *n = h->newNode();
4      n->key = key;
5      n->val = val;
6      setCoords(n, key);
7      for (int i=0; i<D; i++) n->child[i].store(NIL);
8      while (true) {
9          Node* pred = NULL;
10         int dp = 0, dc = 0;
11         s->head = (HeadNode*) (head.load());
12         Node* curr = s->head;
13         LocatePlace(h, dp, dc, pred, curr, n, s);
14         if (dc == D) {
15             // this key is already in the queue
16             return false;
17         }
18         finishInserting(curr, dp, dc);
19         FillNewNode(h, n, dp, dc, curr);
20         uintptr_t temp = (uintptr_t) curr;
21         if (pred->child[dp].compare_exchange_strong(temp, (
22             uintptr_t) n)) {
23             finishInserting(n, dp, dc);
24             RewindStack(s, n, pred, dp);
25             return true;
26         }
27     }
28
29  inline void LocatePlace(Handler* h, int &dp, int &dc,
30                          Node *&pred, Node *&curr,
31                          Node *n, Stack *s) {
32      while (dc < D) {
33          while (curr != NULL && n->k[dc] > curr->k[dc]) {
34              pred = curr;
35              dp = dc;
36              finishInserting(curr, dc, dc);
37              curr = (Node*) (Clear(curr->child[dc].load()));
38          }
39          if (curr == NULL || n->k[dc] < curr->k[dc]) {
40              break;
41          }
42          s->node[dc] = curr;
43          dc++;
44      }
45  }
46
47  inline void FillNewNode(Handler* h, int dp, int dc,
48                          Node* n, Node* curr) {
49      if (dp < dc) {
50          AdoptDesc* desc = h->newDesc();
51          desc->curr = curr;
52          desc->dc = dc;
53          desc->dp = dp;
54          n->adesc.store(desc);
55      } else {
56          n->adesc.store(NULL);
57      }
58      for (int i = 0; i < dp; i++) n->child[i] = F_ADP;
59      for (int i = dp; i < D; i++) n->child[i] = NIL;
60      n->child[dc] = (uintptr_t) curr;
61  }

```

Algorithm 8: Priority Queue Handler

```

1  class Handler{
2      Node* newNode();
3      AdoptDesc* newDesc();
4      Stack* newStack();
5      HeadNode* newHeadNode();
6  };

```

Algorithm 9: Finish Inserting

```

1  void finishInserting(Node *n, int dp, int dc){
2      if (n == NULL) return;
3      AdoptDesc* ad = n->adesc;
4      if (ad == NULL || dc < ad->dp || dp > ad->dc) return;
5      uintptr_t child;
6      Node* curr = ad->curr;
7      for (int i = ad->dp; i < ad->dc; i++) {
8          child = Clear(curr->child[i].fetch_or(F_ADP));
9          uintptr_t temp = NIL;
10         n->child[i].compare_exchange_strong(temp, child);
11     }
12     n->adesc = NULL;
13 }

```

Algorithm 10: Deleting Minimal Node

```

1  uintptr_t deleteMin(Handler* h){
2      Stack* sOld = stack.load();
3      Stack* s = h->newStack();
4      *s = *sOld;
5      int d = D-1;
6      while (true) {
7          Node* last = s->node[d];
8          finishInserting(last, d, d);
9          Node* child = (Node*) (Clear(last->child[d].load()));
10         ;
11         if (child == NULL) {
12             if (d == 0) return NIL;
13             d--;
14             continue;
15         }
16         uintptr_t val = child->val;
17         if (IsMarked(val, F_DEL)) {
18             if (Clear(val) == NIL) {
19                 for (int i = d; i < D; i++)
20                     s->node[i] = child;
21             } else {
22                 s->head = (HeadNode*) (Clear(val));
23                 for (int i = 0; i < D; i++)
24                     s->node[i] = s->head;
25             }
26             d = D-1;
27         } else {
28             if (child->val.compare_exchange_strong(val, F_DEL))
29                 {
30                     for (int i = d; i < D; i++)
31                         s->node[i] = child;
32                     stack.compare_exchange_strong(sOld, s);
33                     int marked = nMarkedNodes.fetch_add(1);
34                     if (marked > R)
35                         purge(s->head, s->node[D-1], h);
36                     return val;
37                 }
38         }
39     }
40 }

```

### C. Details of purge operation

Algorithm 11 outlines the purge operation. Given the head node  $hn$  and the last node to purge  $prg$ , the **purge** operation proceeds only if no other **purge** operation is ongoing (lines 4-5) and if  $hn$  corresponds to the queue's current head node (lines 6-9). The **purge** operation introduces new sentinel head node  $hnNew$  and a copy of  $prg$ ,  $prgNew$ . For each dimension  $d$ , the *LocatePivot* function determines the last node ( $pvt$ ) to be purged at this dimension. If  $pvt.child[d]$  is marked with  $F\_ADP$ , the purge is restarted (lines 23-27). Otherwise,  $pvt.child[d]$  is marked with the  $F\_PRG$  flag (line 57), to prevent it from being changed, and the reference is adapted by either  $hnNew$  or  $prgNew$  (lines 28-38). When all dimensions

Algorithm 11: Purge

```

1 void purge(HeadNode *hn, Node *prg, Handler* h) {
2   if (!notPurging.load()) return;
3   bool temp = true;
4   if (!notPurging.compare_exchange_strong(temp, false))
5     return;
6   if ((uintptr_t) (hn) != head.load()) {
7     notPurging.store(true);
8     return;
9   }
10  nMarkedNodes.store(0);
11  HeadNode* hnNew = h->newHeadNode();
12  Node* prgNew = h->newNode();
13  prgNew->setFromNode(prg);
14  hnNew->val = F_DEL;
15  hnNew->ver = hn->ver + 1;
16  hnNew->key = hn->key;
17  setCoords(hnNew, 0);
18  for (int i=0; i<D; i++) hnNew->child[i].store(NIL);
19  int d = 0;
20  Node* pvt = hn;
21  uintptr_t child;
22  while(d < D) {
23    if (!LocatePivot(prg, pvt, d, child)) {
24      pvt = hn;
25      d = 0;
26      continue;
27    }
28    if (hn == pvt) {
29      hnNew->child[d].store(child);
30      prgNew->child[d].store(F_ADP);
31    } else {
32      prgNew->child[d].store(child);
33      if (d == 0 || prgNew->child[d-1].load() == F_ADP)
34        hnNew->child[d].store((uintptr_t) prgNew);
35      } else {
36        hnNew->child[d].store(NIL);
37      }
38    }
39    d++;
40  }
41  hn->val.store(SetMark((uintptr_t) prg, F_DEL));
42  prg->val.store(SetMark((uintptr_t) hnNew, F_DEL));
43  head.store((uintptr_t) hnNew);
44  Stack* s = h->newStack();
45  UpdateStackAfterPurge(s, hnNew);
46  notPurging.store(true);
47  return;
48 }

50 inline bool LocatePivot(Node* prg, Node* &pvt, int d,
51   uintptr_t &child) {
52   while (pvt->k[d] < prg->k[d]) {
53     finishInserting(pvt, d, d);
54     pvt = (Node*) (Clear(pvt->child[d]));
55   }
56   do {
57     child = pvt->child[d];
58   } while (!IsMarked(child, F_ALL) && !pvt->child[d].
59     compare_exchange_weak(child, SetMark(child, F_PRG)));
60   if (IsMarked(child, F_ADP)) {
61     return false;
62   } else {
63     child = ClearMark(child, F_PRG);
64     return true;
65   }
66 }

```

Algorithm 12: Rewinding deletion stack after insert

```

1 inline void RewindStack(Stack* s, Node* n, Node* pred,
2   int dp) {
3   //NOTE: no need to rewind stack if node is already
4   //deleted...
5   for (bool first_iteration = true; !IsMarked(n->val,
6     F_DEL); first_iteration = false) {
7     Stack* sNow = stack.load();
8     if (s->head->ver == sNow->head->ver) {
9       if (n->key > sNow->node[D-1]->key) {
10        if (!first_iteration) break;
11        *s = sNow;
12      } else {
13        for (int i=dp; i<D; i++) s->node[i] = pred;
14      }
15    } else if (s->head->ver > sNow->head->ver) {
16      Node* prg = (Node*) (ClearMark(sNow->head->val,
17        F_DEL));
18      if (prg->key < sNow->node[D-1]->key) {
19        s->head = (HeadNode*) (ClearMark(prg->val, F_DEL));
20        for (size_t i=0; i<D; i++) s->node[i] = s->head;
21      } else {
22        if (!first_iteration) break;
23        *s = sNow;
24      }
25    } else { // s->head->ver < sNow->head->ver
26      Node* prg = (Node*) (ClearMark(s->head->val, F_DEL));
27      if (prg->key > n->key) {
28        for (int i=dp; i<D; i++) s->node[i] = pred;
29      } else {
30        s->head = (HeadNode*) (ClearMark(prg->val, F_DEL));
31        for (int i=0; i<D; i++) s->node[i] = s->head;
32      }
33    }
34  }
35  if (stack.compare_exchange_strong(sNow, s)) {
36    break;
37  }
38 }

```

have been processed, *hn.val* and *prg.val* are updated with references to help maintaining the deletion stack (lines 41-42) and the deletion stack is updated if needed by function *UpdateStackAfterPurge*.

#### D. Updating deletion stack

After an insert or a purge operation the deletion stack may need to be updated. Algorithm 12 performs such update after inserting a node. Several cases are possible after insert.

**Case 1. The versions of the queue's current stack and the stack after the insertion are the same.**

**Case 1a.** The insert point has lower priority than the last node to be logically deleted. In this case the stack should not be rewound. However, if the queue's stack is older than this insert operation, another concurrent operation, which hasn't seen the effect of this insert, can update the stack and make the inserted node inaccessible. Thus, the stack should be renewed (lines 7-8).

**Case 1b.** The insert point has higher priority than the last node to be logically deleted. In this case the stack must be rewound (line 10).

**Case 2. The version of the queue's current stack is older than the version of the stack after the insertion.**

**Case 2a.** The last node to be logically deleted has lower priority than the *prg* node that corresponds to the stack's



### Algorithm 13: Updating deletion stack after purge

```

1  inline void UpdateStackAfterPurge(Stack* s, HeadNode*
   hnNew) {
2      for (bool first_iteration = true; true;
           first_iteration = false) {
3          Stack* sNow = stack.load();
4          if (hnNew->ver ≤ sNow->head->ver) {
5              // The stack has been updated already
6              return;
7          }
8
9          Node* prg = (Node*) (ClearMark(sNow->head->val,
                                   F_DEL));
10         if (prg->key ≤ sNow->node[D-1]->key) {
11             s->head = (HeadNode*) (ClearMark(prg->val, F_DEL));
12             ;
13             for (size_t i=0; i<D; i++) s->node[i] = s->head;
14         } else {
15             if (!first_iteration) break;
16             *s = *sNow;
17         }
18         if (stack.compare_exchange_strong(sNow, s)) {
19             break;
20         }
21     }
}

```

current version. In this case the stack must be rewound to the next version of the head, which is stored in *prg.val*. It may not be the latest head, but the stack will eventually reach the latest head and every non-deleted node is guaranteed to be accessible (line 16).

**Case 2b.** The last node to be logically deleted has higher priority than than the *prg* node that corresponds to the stack’s current version. In this case the stack should not be rewound but must be renewed as in Case 1a (lines 18-19).

**Case 3. The version of the queue’s current stack is newer than the version of the stack after the insertion.**

**Case 3a.** The item was inserted into purged region (that is, *prg.key* > *n.key*, where *prg* corresponds to the head of the stack after the insertion). In this case the stack must be rewound. It will eventually reach the latest head and every non-deleted node is guaranteed to be accessible (line 24).

**Case 3b.** The item was inserted after *prg*. In this case the stack must be updated to the next version of the head, which is stored in *prg.val*. It may not be the latest head, but the stack will eventually reach the latest node and every non-deleted node is guaranteed to be accessible (lines 26-27).

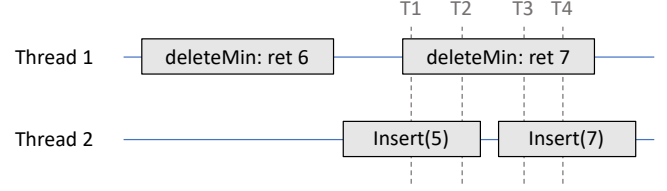
Algorithm 13 demonstrates the procedure of updating the stack after the purge. Because only a *purge* operation can increase the version of the head, the number of possible cases in this situation is smaller than after insertion. The current queue’s stack should only be updated if it is passed the *prg* node that corresponds to the stack’s head (lines 11-12), as in Case 2a for the insertion. If the queue’s current stack is before the *prg* node that corresponds to its version, the stack must be renewed (lines 14-15), as in Case 2b for the insertion.

## V. EVALUATION

### A. Correctness

**Quiescent consistency:** Quiescent consistency is a weaker consistency property than linearizability. According to [4],

(a) Real-time execution sequence



(b) State between  $T3$  and  $T4$

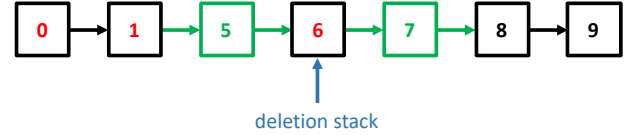


Fig. 1: Example of non-linearizable quiescent consistency.

quiescent consistency requires that method calls separated by a period of quiescence should appear to take effect in their real-time order, and linearizability implies quiescent consistency.

In the implementation presented here, overlapping *insert* operations are linearizable provided that they do not overlap with *deleteMin* operations. A successful CAS operation on line 21 of Algorithm 7 is the linearization point of an *insert* operation with respect to another *insert* operation if the element was inserted. If the element was not inserted because the key was already in the queue, the linearization point is the last load operation on line 37 before the method returns.

Overlapping *deleteMin* operations are also linearizable if they do not overlap with *insert* operations. The linearization point of a *deleteMin* operation with respect to another *deleteMin* operation is either a successful CAS on line 27 of Algorithm 10 if the queue was not empty. If, however, the queue was empty, the linearization point is the last load on line 9 before returning on line 11.

Overlapping *deleteMin* and *insert* methods are quiescently consistent provided that *insert* operations do not overlap with *purge* invocations. An example that demonstrates non-linearizable execution is demonstrated in Figure 1. Provided that the second *deleteMin* operation reads the deletion stack at time  $T1$  while *insert(5)* updates the stack at time  $T2$ , and *insert(7)* adds the item 7 at time  $T3$ , the *deleteMin* method may extract item 7 at time  $T4$  and return this item even though linearizability (as well as sequential consistency [4]) requires that *insert(5)* precedes *insert(7)*.

Even quiescent consistency can be violated if an *insert* operation overlaps with a *purge* operation. Figure 2 demonstrates an example of violation of quiescent consistency. In this example *insert(5)* overlaps with *deleteMin* that initiates a *purge* operation. Assume that the *deleteMin* method extracts item 7 at time  $T1$ , while *insert(5)* inserts item 5 at time  $T2$  before the *purge* operation updates the queue’s head at time  $T3$ . Another instance of *insert*, *insert(2)*, is called after a period

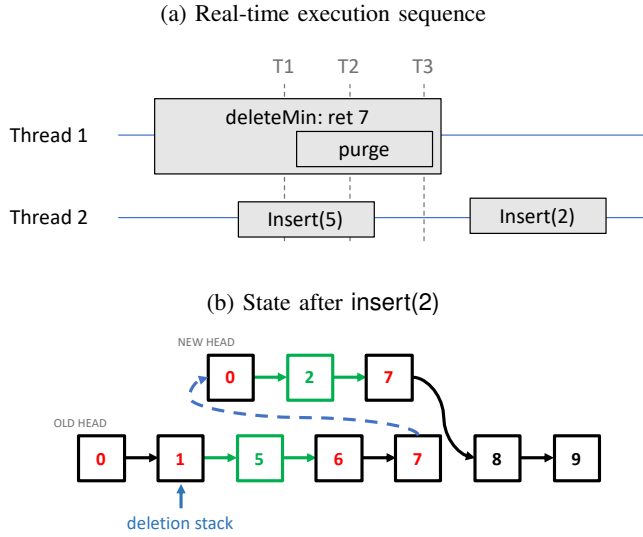


Fig. 2: Example of violation of quiescent consistency.

of quiescence. The method will observe the new head and insert item 2 between 0 and 7. The ensuing state of the queue is shown in Fig. 2b. As a result, two subsequent `deleteMin` invocations, even if separated by periods of quiescence, will first return 5 and only then 2.

## B. Tests

The correctness of the implementation was evaluated by a series of tests, briefly described below, which are essentially the same as those used for the lock-based implementation [11].

*Test 1. Sequential execution:* The first set of tests evaluates correctness of the implementation during sequential execution. In the first stage, elements with all possible priorities in the range from 0 to 262144 except 100 were inserted into the queue, in ascending order. Then an element with priority 100 was inserted. Then `deleteMin` was called repeatedly and the correct order of elements extracted was checked.

The second stage checks whether elements are correctly extracted in ascending order after they are inserted in descending order.

The third stage ensures that the elements are correctly extracted if they have been inserted in a pseudo-random order.

*Tests 2-4. Concurrent execution:* The second test checks that if elements are inserted concurrently by 4 threads, they will be later extracted in the correct order.

The third test ensures that if elements are inserted concurrently and also extracted concurrently afterwards then each extracting thread will observe consistent order of extraction.

The last test checks whether the priority queue operates normally in the case of inserts and extractions by different threads. Each thread during the test performs a randomized mixture of inserts and extractions.

## C. Benchmarks

Performance of the priority queue was evaluated in three mini-benchmarks representing different patterns and ratios

of invocations of `Insert` and `DeleteMin`. All threads in a benchmark are executing similar jobs, which are described in detail in [11].

In the first benchmark, all threads first only insert elements into the queue, then `DeleteMin` operators are called after all threads have finished all insertions.

In the second and third benchmarks, each thread randomly calls `Insert` and `DeleteMin` operations with a given pair of probabilities summing to one. The probability of calling `DeleteMin` is 0.5 for benchmark 2 and 0.2 for benchmark 3. However, to prevent depletion of the queue, if a number of `DeleteMin` operations performed by a thread reaches the number of insert operations performed by the same thread, the thread will perform an `Insert` operation, regardless of a probabilistic draw.

Figure 3 shows the combined performance of all threads (measured in million operations per second – MOPS) of the benchmarks on two different machines and three operating systems (Windows 10 and Ubuntu 18.04.3 LTS on AMD FX-8300 and macOS 10.15 on Intel i5-8259U) when the number of threads is varied from 1 to 32. The performance characteristics were computed based on 10 observations (error bars on the figure indicate standard deviation). On the Windows system, the performance is slightly increasing when the number of threads increases from 1 to 8 and then slightly decreasing with further increase in the number of threads (but the change in the performance does not appear to be substantial).

On the Linux and macOS systems, the performance is increasing for benchmark 1 when the number of threads increases. For benchmarks 2 and 3 on these two systems, the performance stays virtually the same for the number of threads in range from 1 to 4 and in range from 8 to 32. However, it drops significantly when the number of threads increases from 4 to 8.

## VI. CONCLUSIONS AND FUTURE WORK

Overall, in this lock-free version, the performance is about the same when the number of threads increases. We did however see increases in performance in some cases with an increasing number of threads. In contrast, the performance of the lock-based version [11] in this series of papers decreased quickly with increasing numbers of threads.

The presented implementation correctly performs the operations of a priority queue in a multithreaded environment. The performance stays about the same when the number of threads increases. A quantitative comparison between this version and the lock-based implementation in the previous paper in this series [11] will be presented in a future paper.

## REFERENCES

- [1] D. Zhang and D. Dechev, "A Lock-Free Priority Queue Design Based on Multi-Dimensional Linked Lists," *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, no. 3, pp. 613–626, 2016.
- [2] W. Pugh, "Skip lists: A probabilistic alternative to balanced trees," in *WADS 1989: Algorithms and Data Structures*, pp. 437–449, Springer, Berlin, Heidelberg, 1989.



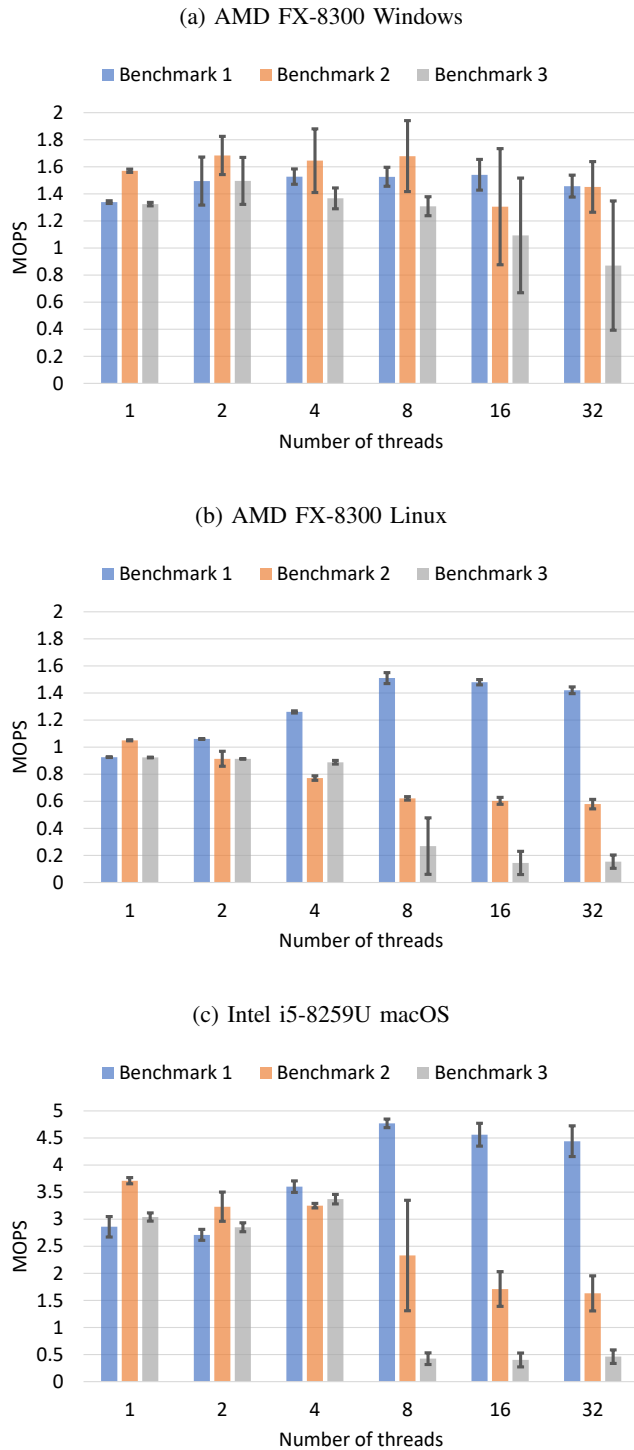


Fig. 3: Average performance (in million operations per second – MOPS) vs number of threads for benchmarks 1-3 executed on three operating systems on two different machines: (a) AMD FX-8300 (8 core 4 units) under Windows 10, (b) same AMD FX-8300 under Ubuntu 18.04.3 LTS and (b) Intel i5-8259U (4 cores 2-way hyper-threading) under macOS 10.15. The results were computed based on 10 observations. Error bars indicate standard deviation.

- [3] H. Sundell and P. Tsigas, “Fast and lock-free concurrent priority queues for multi-thread systems,” *Journal of Parallel and Distributed Computing*, vol. 65, no. 5, pp. 609–627, 2005.
- [4] M. Herlihy and N. Shavit, *The Art of Multiprocessor Programming, Revised Reprint*. Morgan Kaufmann, 2012.
- [5] J. Lindén and B. Jonsson, “A Skiplist-Based Concurrent Priority Queue with Minimal Memory Contention,” in *OPODIS 2013: Principles of Distributed Systems*, pp. 206–220, Springer, Cham, dec 2013.
- [6] F. Ellen, D. Hendler, and N. Shavit, “On the Inherent Sequentiality of Concurrent Objects,” *SIAM Journal on Computing*, vol. 41, pp. 519–536, jan 2012.
- [7] D. Alistarh, J. Kopinsky, J. Li, and N. Shavit, “The SprayList: A scalable relaxed priority queue,” *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP*, vol. 2015-Janua, pp. 11–20, 2015.
- [8] M. Wimmer, J. Gruber, J. L. Träff, and P. Tsigas, “The lock-free k-LSM relaxed priority queue,” *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP*, vol. 2015-Janua, pp. 277–278, 2015.
- [9] H. Rihani, P. Sanders, and R. Dementiev, “MultiQueues: Simpler, Faster, and Better Relaxed Concurrent Priority Queues,” *arXiv*, nov 2014.
- [10] T. L. Harris, “A pragmatic implementation of non-blocking linked-lists,” in *Distributed Computing* (J. Welch, ed.), (Berlin, Heidelberg), pp. 300–314, Springer Berlin Heidelberg, 2001.
- [11] S. Carroll and A. Goopenko, “A C++ implementation of a threadsafe priority queue based on multi-dimensional linked lists and MRLock,” *ResearchGate*, 10 2019.