# Concurrent Priority Queue Using Lock-based Skiplist

Mingbin Xu

Department of Computer Science and Engineering, York University 4700 Keele Street, Toronto, Ontario, Canada, M3J 1P3

**Abstract.** This paper provides an approach of implementing concurrent priority queue using a lock-based skiplist. Skiplist is a probabilistic alternative to balanced trees. Its distributed nature makes it extremely suitable for concurrency. Previous works attempted to wrap a non-blocking concurrent skiplist into a priority queue. However, modification of non-blocking algorithms is generally complex. In this paper, we will review an optimistic lock-based concurrent skiplist implementation, adapt it for concurrent priority queue and evaluate its performance.

## 1 Introduction

The proliferation of multi-core computers have been pervasively affecting how we develop software. It poses increasingly more programming challenges. Data structures and algorithms, while classical, are no longer practical for parallel architectures. Priority queue is one of the fundamental data structures with rich applications, ranging from bandwidth management to discrete event emulation.

Many efforts have been devoted to handling overlapping queue operations in a concurrent setting. Known solutions include a tree-based structure by Shavit and Zemach [9], a heap-based structure by Hunt et al. [5] and recently a skiplist-based structure by Shavit and Lotan [8]. One drawback of the first two approaches is that the underlying data structures require complex rebalancing and centralized lockings, which may cause bottlenecks and contention. As shown in [8], skiplist outperforms other alternatives with relatively simpler implementation. However, the underlying skiplist algorithm credited to Pugh [7] is still considered complicated compared with the one proposed by Herlihy et al. [4] whose synchronization is optimistic and lazy. This paper aims to examine concurrent priority queues based on this simpler synchronization.

## 2 Priority Queue

A priority queue is a set of items where each item is assigned a priority indicating its importance, conventionally the smaller the better. Item priorities are not necessarily distinct and may be time varying. Items with a higher priority are removed before those with a lower priority, irrespective of when they were added. A priority queue must at least support two operations: the *insert* operation to add an item to the set, and the *deleteMin* operation to remove and return the item with highest priority or smallest score. Heap structures, leftist heap [1] and skewed heap [10] for instance, are the most popular backbone for priority queues, giving  $O(\log(n))$  time complexity for both queue operations. Variants such as binomial heap [11] and Fibonacci heap [2] perform faster on certain operations. Besides, priority queues can be modeled by simple linear data structures such as sorted arrays and sorted lists.

Lis	isting 1: PriorityQueue Interface	
1	public interface PriorityQueue <k e<="" th=""><th><pre>xtends Comparable<k>, V&gt; {</k></pre></th></k>	<pre>xtends Comparable<k>, V&gt; {</k></pre>
2	<pre>/** Associates the specified ite</pre>	m(v) with the specified
3	<pre>* priority(k) in this priority</pre>	queue.
4	<pre>* @param k The priority as</pre>	sociated with the new
<b>5</b>	* item.	
6	* @param v The value of th	e item to be stored.
7	* @return The return is imple	mentation-specific. If
8	* duplicate is allowe	d, this method should
9	<ul> <li>always return true;</li> </ul>	otherwise, when
10	<ul> <li>* duplicate is hit, i</li> </ul>	t does NOT put the item
11	* in the queue and re	turns false. */
12	public boolean insert(K k, V v	);
13		
14	<pre>/** Deletes and returns an item</pre>	with highest priority.
15	* ©return An item with hi	ghest priority. */
16	<pre>public V deleteMin();</pre>	
17	}	

### 3 SkipList



Fig. 1: This example is a skiplist of 6 items. The bottom list contains all items. The lists above are the sublists of the lists below them. In order to search for the item 8, traversal begins at the leftmost node of the topmost level. At each level, move rightwards until a key greater than or equal to what is searched for is hit. If such a key has been hit, go one layer down and repeat this process. As seen in this example, an "express" from 0 to 7 at the second topmost level skips 2 and 5.

Skiplist [6] is a probabilistic data structure that allows efficient search, insertion and removal. A skiplist maintains a collection of sorted linked hierarchy, which

#### $\mathbf{2}$

mimics, with some randomization, a balanced search tree. Searching starts with the sparsest linked list in a top-down fashion, traversing until the largest item smaller than the sought is found at each level. Fast searching is made possible by skipping over a few elements at each level. Figure 1 shows a skiplist with integer keys. The higher-level lists serve as expresses to the lower-level lists. The bottom level is a regular sorted linked list with all items. Each item has a probability of 0.5 to be chosen in the level above. More formally, an item is present at level i with probability  $0.5^i$ . For example, the nodes at level 0 have a chance of 0.5 to appear at level 1, and each node at level 1 will be promoted to level 2 with probability  $0.5^2$ . Thus, roughly speaking, level i contains  $n * 0.5^i$  items, and each link at level i skips  $2^i$  items, where n is the number of items in the list.

The structures of a skiplist and a skiplist node are detailed in Figure 3 (UML) and Listing 6 (Appendix). It is convenient to add two sentinels \_\_head and \_\_tail with priorities  $-\infty$  and  $+\infty$  respectively. Each skiplist node, be it sequential or concurrent, stores an array of skiplist node references which are its successors at all levels. Listing 2 shows the \_\_getPredecessorsAndSuccessors method. It traverses the skiplist with *predecessors* and *successors* arrays starting from *\_\_head* at the highest level. Given target to be searched, the \_\_getPredecessorsAndSuccessors method descends from one level to another (Line 15-26). At each level it continuously sets *current* to be the successor of *predecessor* until it encounters a node whose priority is larger than the one sought (Line 17-20). If it locates a node with a matching priority, it records the height of the node in *levelFound* (Line 21-23). Before going one level down, \_\_getPredecessorsAndSuccessors also puts predecessor and successor in the arrays of predecessors and successors (Line 24-25). If it finds target, predecessor itself is target; otherwise it is the largest node that is smaller than *target*. Figure 1 demonstrates the process of *search*. When performing *insert* and *remove*, one could keep track of the predecessors and the successors of the search path and then re-connect the predecessors and the successors, as depicted in Figure 2.

It is straightforward to adapt skiplist to implement a priority queue of items tagged with priorities. Skiplist is a collection of ordered items, ensuring that high-priority items appear at the front of the list. The insertion of a skiplist-based priority queue therefore reuses the interface from a skiplist without any modification. Since the bottom level contains all items, the very next node to which *\_\_head* points, if not *\_\_tail*, is the smallest.

Listing 2:getPredecessorsAndSuccessors of $LazySkipListPriorityQueue$ and $CoarseSkipListPriorityQueue$				
1	/**	<pre>@param</pre>	target The node to be sought.	
2	*	<pre>@param</pre>	predecessors Used to retrieved the largest nodes at each	h
3	*		level NOT greater than target.	
4	*	<pre>@param</pre>	successors Used to retrieved the smallest nodes at eac	ch
<b>5</b>	*		level greater than target.	
6	*	<b>@return</b>	It guarantees to return -1 if invoked by insert() because	
7	*		Node class imposes uniqueness; between 0 and 32 (inclusive)	)
8	*		if invoked by deleteMin() because deleteMin() removes the	
9	*		very node which must exist. */	
10	priv	vate int	getPredecessorsAndSuccessors (	

```
Node<K, V> target, Node<K, V>[] predecessors,
11
             Node <K, V>[] successors ) {
^{12}
        int levelFound = -1;
13
        Node <K, V> predecessor =
14
                                    __head;
        for ( int i = MAX_LEVEL; i >= 0; i-- ) {
15
             Node<K, V> current = predecessor.__next[i];
16
             while ( current != __tail && target.compareTo(current) > 0 ) {
17
^{18}
                 predecessor = current;
                 current = predecessor.__next[i];
^{19}
20
             }
                ( levelFound == -1 && target.compareTo(current) == 0 ) {
             if
^{21}
                 levelFound = i;
^{22}
             }
23
             predecessors[i] = predecessor;
^{24}
             successors[i] = current;
^{25}
        }
^{26}
        return levelFound;
27
        // end of __getPredecessorsAndSuccessors
    }
^{28}
```



(a) Suppose that 6 is going to be inserted. The search algorithm gives a search path highlighted in red. The starting/finishing points of the red arrows are the predecessors/successors of the new item. Insert the new node in-between.



(b) Suppose that 6 is going to be removed. Similarly, the search algorithm returns a group of predecessors (starting point of the red arrows). Relink these predecessors' successors to 6's successors (the finishing points of outgoing arrows from 6).

Fig. 2: Insertion and Removal

# 4 Concurrent Skiplist-Based Priority Queue and Implementation

We now turn our attention to skiplist's concurrent implementation in Java. The same trick from [3] is applied to skiplist, that is, the *deleteMin* operation is done lazily with two steps: first label the target node, removing it logically, and second, reconnect its predecessors' next fields, removing it physically. Each node contains two *Boolean* fields, *\_\_logicallyDeleted* and *\_\_fullyConnected*. The former is to indicate that the node has been logically deleted and is ready to be physically deleted. The latter is to notify that a node is being inserted but not



Fig. 3: Class Relationship in UML. The Java Implementations of *LazySkipListPriorityQueue* is given in the Appendix.

yet well-connected, and therefore is not considered in the list. The node to which \_\_head points is not necessarily smallest since it might not be *fullyConnceted* or be marked as *logicallyDeleted*.

The algorithm is optimistic. Similar to [3], searching for the insertion spot or the victim to be removed does not acquire locks. During *insert* and *deleteMin*, an item is spliced in or out only when its predecessors are fully connected and not logically deleted, and their successors remain unchanged; it retries otherwise. Predecessors are retrieved by a top-down search. However, *deleteMin* scans the bottom level from the list head. If duplicated priority is allowed, the former finds the highest match. Contrarily, the latter stops as soon as it hits the first valid match. They are not necessarily the same. Uniqueness is enforced by assigning a creation time to each node upon its construction so as to avoid inconsistent search.

The insert method, shown in Listing 3, initializes the predecessors and the successor arrays before calling  $\_\_getPredecessorsAndSuccessors$  (Line 5-6).  $\_\_getPredecessorsAndSuccessors$  determines the predecessors and the successors of a new node and put them into the said arrays respectively. These two arrays are not reliable because they may not be accurate when the nodes are accessed. The thread executing insert acquires locks in descending priority (bottom-up) and validates each of the predecessors (Line 18-30). The validation at each level checks whether predecessor is still adjacent to successor and neither is marked as logicallyDeleted. If validation fails, there must be some contention introduced by other threads. The current thread releases the locks it holds and simply retries. If the current thread manage to locks all the predecessors, it splices in target between predecessors and successors (Line 35-38). When all re-connections are done, the field  $\_\_fullyConnected$  of target is set to true, which indicates that target is inserted in the queue (Line 41).

```
Listing 3: insert of LazyConcurrentSkipListPriorityQueue
    /** {@inheritDoc} */
   public boolean insert(K key, V value) {
2
        int newHeight = __randomLevel();
з
4
        Node <K, V> target = new Node <K, V>(newHeight, key, value);
        Node<K, V>[] predecessors = (Node<K, V>[]) new Node[MAX_LEVEL + 1];
5
        Node <K, V>[] successors = (Node <K, V>[]) new Node [MAX_LEVEL + 1];
6
7
        while ( true ) {
8
            int levelFound = __getPredecessorsAndSuccessors(
9
10
                                     target, predecessors, successors );
            assert levelFound == -1;// Uniqueness is enforced by Node class.
11
12
            // keep track of the lock height, so as to know where to unlock
13
            int lockHeight = -1;
14
15
16
            try {
                 boolean valid = true;
17
                Node<K, V> predecessor, successor;
18
                for ( int i = 0; valid && i <= newHeight; i++ ) {</pre>
19
                    predecessor = predecessors[i];
20
                    successor = successors[i];
21
                    predecessor.__lock.lock();
22
                    lockHeight = i;
23
```

```
^{24}
                      // 1. If predecessor.__logicallyDeleted, predecessor
^{25}
                         is being removed.
^{26}
                      // 2. If successor.__logicallyDeleted, successor is
27
                            being removed.
^{28}
                      11
                      // 3. If predecessor.__next[i] != successor, sth has
29
                           been inserted in-between
                      11
30
31
                      // In either case, the predecessor or successor is invalid.
                      valid = !predecessor.__logicallyDeleted &&
32
33
                               !successor.__logicallyDeleted &&
                               predecessor.__next[i] == successor;
34
                 }
35
                 if
                    ( !valid ) {
36
                      continue;
37
                 }
38
39
                 for ( int i = 0; i <= newHeight; i++ ) {</pre>
40
                     target.__next[i] = successors[i];
41
                     predecessors[i].__next[i] = target;
42
^{43}
^{44}
                 // linearization point, the node is considered "inserted".
45
                 target.__fullyConnected = true;
46
                 return true:
47
             }
                // end of try
48
             finally {
49
                 for ( int i = 0; i <= lockHeight; i++ ) {</pre>
50
                     predecessors[i].__lock.unlock();
51
52
                 }
53
             }
                 // end of finally
        }
             // end of while ( true )
54
   }
        // end of insert
55
```

One can disallow deleteMin to remove a node whose creation time is earlier than the function invocation if the FIFO (first in first out) convention of a queue needs to be preserved. More formally, let dm be a deleteMin operation, Q be the item set whose elements invoke insert before dm, D be the item set that is removed by other deleteMin operations that precede or are concurrent with dm, and I be the item set that is concurrently inserted during dm's execution. If timestamps are preserved, dm returns the minimum from Q - D; Q + I - Dotherwise.

The deleteMin method appears in Listing 4. At the beginning (Line 3), it records the invocation time, nodes of creation time greater than which will not be considered. Similarly two arrays are initialized to retrieve the predecessors and the successors (Line 5-6). deleteMin scans the bottom list to get the smallest item (Lines 11-24). If current is ready to be deleted, deleteMin acquires lock, sets the field \_\_logicallyDeleted to true (Line 14-20); otherwise current becomes the predecessor of current (Line 22-23). Note that new arrivals may be inserted before current, but their creation time cannot be earlier than deleteMin's invocation time. If \_\_tail is encountered, the queue is empty before the start of deleteMin, though it might not be the case during the execution of deleteMin. \_\_getPredecessorsAndSuccessors is called to retrieve target's ostensible predecessors and successors. In order to avoid deadlock, locks are acquired and released in the same manner as insert (bottom-up, Line 38-44, 54-58). predecessor is verified as valid when it points to *target* and is not logically deleted (Line 42-43). Once all levels are locked and valid, *target* is spliced out (Line 45-50).

A small modification is made in Listing 5 where time restriction on *deleteMin* is not imposed. When scanning the bottom list, the *while* loop is broken by *timeInvoked* in Line 13 Listing 4 but is not in Line 8 Listing 5. Another change is in Line 17 Listing 5 where search is retried from the head of the list. We will see in the experiment section that this impacts the throughput.

```
Listing 4: deleteMin of LazyConcurrentSkipListPriorityQueue
    /** {@inheritDoc} */
 1
    public Node<K,V> __deleteMin() {
    long timeInvoked = TIMER.get();
    int targetHeight = -1;
    ())
2
3
 4
         Node<K, V>[] predecessors = (Node<K, V>[]) new Node[MAX_LEVEL + 1];
Node<K, V>[] successors = (Node<K, V>[]) new Node[MAX_LEVEL + 1];
 \mathbf{5}
 6
         Node<K, V> predecessor, current, target = null;
 7
 8
9
         predecessor = __head;
         current = __head.__next[0];
10
         while ( current != __tail ) {
    if ( current.__fullyConnected && !current.__logicallyDeleted
11
12
                             && current.__timeCreated > timeInvoked ) {
13
14
                   current.__lock.lock();
15
                   if ( !current.\_logicallyDeleted ) {
16
                        current.__logicallyDeleted = false;
17
                        current.__lock.unlock();
^{18}
                        break;
                   3
19
^{20}
                   current.__lock.unlock();
              }
^{21}
22
              predecessor = current;
              current = predecessor.__next[0];
^{23}
^{24}
         }
^{25}
26
         if ( current != __tail ) {
              while ( true ) {
27
                   int levelFound =
^{28}
                   __getPredecessorsAndSuccessors(
^{29}
                                      current, predecessors, successors );
30
                   target = successors[levelFound];
31
                   assert ( levelFound != -1 && current == target );
32
33
^{34}
                   target.__lock.lock();
35
                   int lockHeight = -1;
36
                   try {
                        boolean valid = true;
37
                        for ( int i = 0; valid && i <= targetHeight; i++ ) {</pre>
38
                             predecessor = predecessors[i];
39
                             predecessor.__lock.lock();
40
                             lockHeight = i;
41
                             valid = !predecessor.__logicallyDeleted &&
^{42}
                                           predecessor.__next[i] == target;
^{43}
                        }
^{44}
                        if ( !valid ) {
45
                             continue;
46
                        }
47
                        for ( int i = 0; i <= targetHeight; i++ ) {</pre>
^{48}
                             predecessors[i].__next[i] = target.__next[i];
49
                        }
50
                        target.__lock.unlock();
51
                        return target;
52
                   }
53
                   finally {
54
```

```
55
                     for ( int i = 0; i <= lockHeight; i++ ) {
                          predecessors[i].__lock.unlock();
56
57
                     }
                 }
58
59
             }
                 //
                    end of while ( true )
60
                end of if ( current != __tail )
        }
             11
61
62
        return null;
    }
        // end of deleteMin
63
```

```
\label{eq:listing 5: deleteMin of LazyConcurrentSkipListPriorityQueue
```

```
/** {@inheritDoc} */
 1
    public Node<K,V> __deleteMin() {
    // see listing 4
2
3
         // .....
 4
         predecessor = __head;
 \mathbf{5}
         current = __head.__next[0];
 6
         while ( current != __tail ) {
    if ( current.__fullyConnected && !current.__logicallyDeleted ) {
 7
 8
 9
                   current.__lock.lock();
10
                   if ( !current.__logicallyDeleted ) {
11
                        current.__logicallyDeleted = false;
^{12}
                        current.__lock.unlock();
^{13}
                        break;
                   }
14
15
                   else {
16
                        current.__lock.unlock();
17
                        current = __head;
^{18}
                   }
19
              }
^{20}
              predecessor = current;
^{21}
              current = predecessor.__next[0];
^{22}
         }
^{23}
         //
^{24}
         //
             see listing 4
^{25}
    }
```



(a) Suppose that 6 is going to be inserted. Step 1: Identical to a sequential skiplist, the search algorithm returns the predecessors and the successors of the new node, which is indicated by the red arrows. Step 2: From right to left, each predecessor, if valid (fully connected, not logically deleted, and still pointing to the successor returned by the search algorithm), is locked; otherwise, insertion unlocks these predecessors and retries from Step 1.



(b) Step 3: Splice in the new node from right to left. As depicted here, 6 has been linked between 5 and 7 in Level 0, Level 1 and Level 2. 0 still points to 7. 6 is not fully connected, and thus is not considered in the skiplist.

Fig. 4: Concurrent *insert* 



(a) Suppose 0, 2 and 5 were inserted after *deleteMin*'s invocation and 6 was inserted before *deleteMin*. Step 1: Scan the bottom list from left to right until the first node inserted earlier than *deleteMin*'s invocation and not logically deleted is found. In this example, 6 hasn't been logically deleted. *\_\_remove* sets it as logically deleted (red) and locks it.



(b) Step 2: from right to left, each predecessor, if valid (not logically deleted and still pointing to the successor returned by the search algorithm), is locked; otherwise, *deleteMin* unlocks those predecessors and retries from Step 1.



(c) Step 3: from right to left, link each predecessor to 6's successor at the same level. Unlock the predecessors when all done.

Fig. 5: Concurrent deleteMin

## 5 Experiment

#### 5.1 Experiment Setup

The proposed algorithm is implemented as a class LazySkipListPriorityQueue in the Java programming language. Its performance is evaluated by comparing with a priority queue wrapping ConcurrentSkipListMap class from the *java.util.concurrent* package<sup>1</sup>. In addition, we also provide a baseline implementation in which methods are synchronized to guarantee thread-safety and a dummy priority queue whose methods are synchronized but without any actual implementation. Since deleteMin introduces heavy contention on the head of a list, we run the algorithms with insert-to-deleteMin ratio of 5:5, 7:3 and 9:1 in the priority range from 0 up to 100, 1000 and 10000, and examine how deleteMin affects performance. All experiments starts with empty data structures with priorities selected from a uniform distribution. An operation is decided randomly with bias according to the insert-to-deleteMin ratio. The experiments are conducted in the Manycore Testing Lab  $(MTL)^2$ . It is a remote system equipped with 256GB memory and 4 CPUs of Intel<sup>®</sup> Xeon<sup>®</sup> E7-4860 @ 2.27GHz. In our experiments, we constrained the CPU usage on the first 16 logical processing units.<sup>3</sup> and parallel executions up to 32 threads.

<sup>&</sup>lt;sup>1</sup> We also implemented and compared with heap-based priority queue credited to Hunt et al. [5]. However, as concurrency rose, a single operation sometimes did not finish within 1 second. It failed to provide a valuable reference. Moreover, it became slower than a sequential heap when the thread number exceeded 3.

<sup>&</sup>lt;sup>2</sup> Intel Corporation has set up a special remote system that allows faculty and students to work with computers with lots of cores, called the Manycore Testing Lab. Users intentionally write programs that take advantage of multi-core parallelism and explore the issues in parallelism and concurrency that arise.

<sup>&</sup>lt;sup>3</sup> Intel<sup>®</sup> Xeon<sup>®</sup> E7-4860 each has 10 cores and meanwhile supports hyper-threading. The system therefore is of 80 logical processing units.

#### 5.2 Results & Analysis

Provided that the monitor baseline and the lock-free wrapper do not take into account the temporal relationship between *insert* and *deleteMin*, we would like to investigate the throughput with and without the timestamps assigned to an item and the invocation of *deleteMin* to ensure fairness. Results are shown in Figure 2-4. Graphs in the right columns depict the throughput of various algorithms and the left columns detail the number for retrying and the ratio of the accumulative numbers of nodes scanned at the bottom list by *deleteMin* to the number of *deleteMin*. The following conclusions have been reached<sup>4</sup>:

- The *deleteMin* operation brings in a great deal of contention. As the ratio of *insert* grows, contention is alleviated and thus the throughput increases. It conforms to our expectation that *deleteMin* is centralized around the list head while *insert* is well distributed.
- Overhead is primarily in the form of locking. Retrying, though important to ensure correct concurrency, is a rare case which introduces little overhead. As we can see in the left columns, retrying is outnumbered several thousands to one by operation count.
- Priority range does not give rise to contention, because uniqueness is enforced on the *Node* class by an atomic counter. Experiments are divided into three groups in the right columns. They behave similarly in all three different ratio settings.
- If deleteMin is not allowed to remove items whose creation time is earlier than the invocation of deleteMin, deleteMin scans a significantly longer path possibly leading to the tail of the list. At a high concurrent level, context switch is likely to take place right after deleteMin's invocation and many new arrivals jump the queue. However, these new arrivals are close to the head, thus does not result in a lengthy search path. It explains why LazySkipListPriority with timestamps is of lower throughput.
- LazySkipListPriorityQueue yields comparable but more predictable result to that of the lock-free wrapper. There is no obvious winner when deleteMin dominates (5:5 ratio). Both lock-free and lock-based implementation are close to monitor in such situation.

## 6 Conclusion

This paper explores a friendly implementation of concurrent priority queue. search, insert and remove of concurrent skiplist, in common usage, spread all

<sup>&</sup>lt;sup>4</sup> The fluctuation of the lock-free wrapper is so far not explainable. It might be a hardware problem. A reasonable guess of monitor's performance gain is that the optimizer moves some local statements out of the synchronized block.



Fig. 6: Throughput within 1 second in a range of 100 and the corresponding number of retrying and the number of nodes scanned by deleteMin in the bottom list



Fig. 7: Throughput within 1 second in a range of 1000 and the corresponding number of retrying and the number of nodes scanned by deleteMin in the bottom list



Fig. 8: Throughput within 1 second in a range of 10000 and the corresponding number of retrying and the number of nodes scanned by deleteMin in the bottom list

along the list. Therefore they gain performance as a matter of course. In contrast, a skiplist-based priority queue centralizes its removal on the very first item, which is heavily contentious. If *deleteMin* is not allowed to remove items whose creation time is earlier than the invocation of *deleteMin*, *deleteMin* scans a seriously longer path. Time is spent on an actual traversal instead of concurrent overhead. Both implementations do not benefit much from multi-threads when *insert* and *deleteMin* are equally mixed.

#### Acknowledgement

We thank the team at Intel's Multicore Testing Lab for providing us access to their machine.

#### References

- 1. Clark Allan Crane, Linear lists and priority queues as balanced binary trees, (1972).
- Michael L Fredman and Robert Endre Tarjan, Fibonacci heaps and their uses in improved network optimization algorithms, Journal of the ACM (JACM) 34 (1987), no. 3, 596–615.
- Steve Heller, Maurice Herlihy, Victor Luchangco, Mark Moir, William N Scherer III, and Nir Shavit, A lazy concurrent list-based set algorithm, Principles of Distributed Systems, Springer, 2006, pp. 3–16.
- Maurice Herlihy, Yossi Lev, Victor Luchangco, and Nir Shavit, A simple optimistic skiplist algorithm, Structural Information and Communication Complexity, Springer, 2007, pp. 124–138.
- Galen C Hunt, Maged M Michael, Srinivasan Parthasarathy, and Michael L Scott, An efficient algorithm for concurrent priority queue heaps, Information Processing Letters 60 (1996), no. 3, 151–157.
- William Pugh, Skip lists: a probabilistic alternative to balanced trees, Communications of the ACM 33 (1990), no. 6, 668–676.
- 7. \_\_\_\_\_, Concurrent maintenance of skip lists, (1998).
- Nir Shavit and Itay Lotan, Skiplist-based concurrent priority queues, Parallel and Distributed Processing Symposium, IEEE, 2000, pp. 263–268.
- Nir Shavit and Asaph Zemach, *Diffracting trees*, ACM Transactions on Computer Systems (TOCS) 14 (1996), no. 4, 385–428.
- Daniel Dominic Sleator and Robert Endre Tarjan, Self-adjusting heaps, SIAM Journal on Computing 15 (1986), no. 1, 52–69.
- Jean Vuillemin, A data structure for manipulating priority queues, Communications of the ACM 21 (1978), no. 4, 309–315.

# A Java Implementation of Data Structure Discussed

```
import java.io.*;
1
    import java.util.concurrent.locks.*;
2
   import java.util.concurrent.atomic.*;
import java.util.*;
3
 4
 \mathbf{5}
    /**
 6
    * A class which implements PriorityQueue interface and whose underlying
7
    * data structure is SkipList. This implementation allows duplicated
 8
9
    * priority.
    * @param <K>
                    The type of priority that is associated to an item.
10
   * @param <V>
                    The type of item stored in the priority queue.
11
   * */
12
    public class LazySkipListPriorityQueue <K extends Comparable<K>, V>
13
    implements PriorityQueue<K, V> {
14
15
16
        * Each node is associated with a timestamp. The reason is as follow:
17
        * In the first stage of deleteMin(), it begins traversal at the
18
        * bottom; therefore it finds the very FIRST match. In the second
19
        * stage of deleteMin(), it starts traversal at the top; therefore it
20
        * finds the HIGHEST match, which is NOT necessary. the first match.
21
22
        \ast By adding a sorted timestamp, uniqueness is enforced, i.e., the
23
        * underlying data structure has no duplicate.
        * hashCode() though enforces uniqueness, sort order is not guaranteed.
24
25
        * */
^{26}
        final private static class Node<K, V> implements Comparable<Node<K, V>> {
27
                                                    // the client-provided priority
28
             final private K __key;
             final private V __value;
29
30
             final private int __level;
                                                    // the height of the node
             final private Node<K,V>[] __next;
^{31}
32
             final private Lock __lock;
33
34
             // if true, the node is not considered existed
             // false by default; once changed to true by deleteMin, % \left( {{{\left( {{{\left( {{{\left( {{{}_{{\rm{m}}}}} \right)}} \right)}_{{\rm{m}}}}}} \right)
35
             // it cannot become false again
36
             volatile private boolean __logiciallyDeleted;
37
38
             // if true, the node is considered in the list
39
             // false by default, once changed to true by insert,
40
             // it cannot become false again
41
42
             volatile private boolean __fullyConnected;
43
             // linearized point at which method "takes effect"
44
            final private long __timeCreated; // also serves as a unique id
^{45}
46
             @SuppressWarnings("unchecked")
47
^{48}
             * Create a node whose height is level, priority is key and item
49
50
             * stored is value
             * @param level
                                 Height of the node.
51
             * @param key
                                 Priority.
52
             * @param value
                                 Item to be stored.
53
             * */
54
             public Node( int level, K key, V value ) {
55
                 __level = level;
56
                 __key = key;
57
                 __value = value;
58
                 __timeCreated = TIMER.getAndIncrement();
59
                 __logiciallyDeleted = false;
60
                 __next = (Node<K, V>[]) new Node[level + 1];
61
                 __lock = new ReentrantLock();
62
```

```
__fullyConnected = false;
} // end of Node
 63
 64
 65
 66
                @Override
 67
                /** {@inheritDoc} */
 68
 69
                public int compareTo(Node<K, V> o) {
                     // only the keys of __head and __tail are null
// only __tail could be passed to this function
 70
 71
                     // __tail is a sentinel node which is always the largest if ( o.\__key == null ) {
 72
 73
                          return -1;
 74
                     }
 75
 76
 77
                     @SuppressWarnings("unchecked")
                     int result = ((Comparable<K>)__key).compareTo(o.__key);
if ( result == 0 ){
 78
 79
                          if ( __timeCreated < o.__timeCreated ) {</pre>
 80
                              return -1;
 81
                          }
 82
                          else if ( __timeCreated > o.__timeCreated ) {
   return 1;
 83
 84
                          }
 85
                          else {
 86
                              return 0:
 87
                          }
                              // end of if ( __timeCreated )
 88
                     }
 89
                     else {
 90
 91
                         return result;
 92
                   } // end of if ( result )
// end of compareTo
 93
                }
          }
               // end of class Node
 94
 95
 96
 97
          // The list is sorted by a combination of (key, timeinserted). final private static int MAX_LEVEL = 32;
 98
99
           final private static AtomicInteger TIMER = new AtomicInteger();
100
101
          // 2 sentinel nodes. __head is "smallest". __tail is "largest".
final private Node<K, V> __head;
102
103
           final private Node<K, V> __tail;
104
105
106
107
           /**
108
           * Default Constructor which creates an empty concurrent priority
109
           * queue.
110
           * */
111
           public LazySkipListPriorityQueue() {
               __head = new Node<K, V>(MAX_LEVEL, null, null);
__tail = new Node<K, V>(MAX_LEVEL, null, null);
112
113
                for ( int i = 0; i < __head.__next.length; i++ ) {
    __head.__next[i] = __tail;</pre>
114
115
116
                }
117
           }
                // end of LazySkipList
118
119
           /**
120
121
           * Generates a random integer between 0 and 32, inclusive. The
           * probability that i is returned is 0.5^i if i is not 32 and
122
           * 0.5<sup>31</sup> if i is 32.
123
           * @return A random integer.
124
125
           * */
126
           private static int __randomLevel() {
               int height = 0;
127
                while ( height < MAX_LEVEL && Utility.randomBoolean() ) {</pre>
128
129
                     height++;
                }
130
```

```
131
                 return height;
           } // end of randomHeight
132
133
134
            /** Please look at Listing 2 */
135
           private int __getPredecessorsAndSuccessors ( Node<K, V> target,
Node<K, V>[] predecessors, Node<K, V>[] successors ) {
136
137
138
               // please look at Listing 2
           } // end of __getPredecessorsAndSuccessors
139
140
141
            @SuppressWarnings("unchecked")
142
           @Override
143
            /** {@inheritDoc} */
144
           public boolean insert(K key, V value) {
145
                // Please look at Listing 3
146
           } // end of insert
147
148
149
            @SuppressWarnings("unchecked")
150
151
           @Override
            /** {@inheritDoc} */
152
           public V deleteMin() {
153
           // Please look at Listing 4
} // end of deleteMin
154
155
156
157
158
           /**
159
           * Return the string representation of the current object. This
* method is written for test purpose and is NOT thread-safe.
* @return The string representation of the current object, of
160
161
162
                         which is a sorted sequence. */
163
            public String toString() {
164
                 StringBuffer sb = new StringBuffer();
sb.append("[ ");
for ( Node<K, V> current = __head.__next[0];
165
166
167
                      current != __tail; current = current.__next[0] ) {
sb.append('(').append(current.__key).append(", ")
.append(current.__timeCreated).append(", ")
168
169
170
171
                         .append(current.__value).append(") ");
172
                 }
                 sb.append(']');
173
                 return sb.toString();
174
175
           }
               // end of toString()
176
177
178
179
            public boolean isValid() {
                 TreeSet<K> prevLevel = new TreeSet<K>();
TreeSet<K> curLevel = null;
180
181
182
                 for ( int i = 0; i <= MAX_LEVEL; i++ ) {</pre>
183
184
                      curLevel = new TreeSet <K>();
                      Node<K, V> prev = __head;
Node<K, V> cur = __head.__next[i];
boolean result = true;
185
186
187
188
                      while ( cur != __tail) {
189
                      if ( !cur.__fullyConnected || cur.__logiciallyDeleted ) {
190
191
                           System.err.println(
                                 i + " Fail because of invalid boolean flags..." );
192
193
                           result = false;
194
                      3
195
                      if ( prev != __head && prev.compareTo(cur) >= 0 ) {
                           System.err.println(
196
                                 i + " Fail because of sort order..." );
197
                            result = false;
198
```

```
199
                       }
                       if ( !result ) {
200
201
                             // this piece of code sees how sort order is violated
202
                             // the console has only limited buffer. in order to
                            // observe the result, we have to write it to a file
File file = new File("./logging");
203
204
205
                             BufferedWriter bw = null;
206
                            try {
                                  if ( !file.exists() ) {
207
                                       file.createNewFile();
208
209
                                  }
                                  bw = new BufferedWriter(
210
                                            new FileWriter(file.getAbsoluteFile()));
211
212
                                  Node<K, V> previous = __head;
Node<K, V> current = __head.__next[i];
while ( current != __tail ) {
    bw.write( current.__key + " " +
        current.__timeCreated + "\n" );
    if ( previous != head &&
213
214
215
216
217
                                       if ( previous != __head &&
    previous.compareTo(current) >= 0 ) {
218
219
                                             System.err.println("violate skiplist property");
bw.write("violate skiplist property\n");
220
221
                                       }
222
                                       previous = current;
current = current.__next[i];
223
224
                                  3
225
                                  bw.close():
226
227
                                  System.exit(1);
228
                            }
                            catch (IOException e) {
229
                                  // let it throw...
230
                            }
231
                            return false;
} // end of if ( !result )
232
233
234
                            curLevel.add( cur.__key );
235
                            prev = cur;
cur = cur.__next[i];
236
237
                       7
                            // end of while ( cur != __tail)
238
239
240
                       if ( i > 0 ) {
                             // check if higher level is subset of the lower level
241
                            if ( !prevLevel.containsAll( curLevel ) ) {
    System.err.println( i + " Fail because of xxoo" );
242
^{243}
244
                                  return false;
245
                            }
                      }
246
247
                      prevLevel = curLevel;
^{248}
                 }
                       // end of for
249
                 return true;
                // end of isValid
250
            }
251
252 }
           // end of LazySkipListPriorityQueue
```