Jiffy: A Fast, Memory Efficient, Wait-Free Multi-Producers Single-Consumer Queue

Dolev Adas and Roy Friedman
Computer Science Department
Technion
{sdolevfe,roy}@cs.technion.ac.il

November 3, 2020

Abstract

In applications such as sharded data processing systems, sharded in-memory key-value stores, data flow programming and load sharing applications, multiple concurrent data producers are feeding requests into the same data consumer. This can be naturally realized through concurrent queues, where each consumer pulls its tasks from its dedicated queue. For scalability, wait-free queues are often preferred over lock based structures.

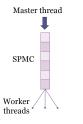
The vast majority of wait-free queue implementations, and even lock-free ones, support the multi-producer multi-consumer model. Yet, this comes at a premium, since implementing wait-free multi-producer multi-consumer queues requires utilizing complex helper data structures. The latter increases the memory consumption of such queues and limits their performance and scalability. Additionally, many such designs employ (hardware) cache unfriendly memory access patterns.

In this work we study the implementation of wait-free multi-producer single-consumer queues. Specifically, we propose Jiffy, an efficient memory frugal novel wait-free multi-producer single-consumer queue and formally prove its correctness. We then compare the performance and memory requirements of Jiffy with other state of the art lock-free and wait-free queues. We show that indeed Jiffy can maintain good performance with up to 128 threads, delivers up to 50% better throughput than the next best construction we compared against, and consumes $\approx 90\%$ less memory.

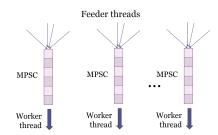
1 Introduction

Concurrent queues are a fundamental data-exchange mechanism in multi-threaded applications. A queue enables one thread to pass a data item to another thread in a decoupled manner, while preserving ordering between operations. The thread inserting a data item is often referred to as the *producer* or *enqueuer* of the data, while the thread that fetches and removes the data item from the queue is often referred to as the *consumer* or *dequeuer* of the data. In particular, queues can be used to pass data from multiple threads to a single thread - known as *multi-producer single-consumer queue* (MPSC), from a single thread to multiple threads - known as *single-producer multi-consumer queue* (SPMC), or from multiple threads to multiple threads - known as *multi-producer multi-consumer queue* (MPMC). SPMC and MPSC queues are demonstrated in Figure 1.

MPSC is useful in sharded software architectures, and in particular for sharded in-memory key-value stores and sharded in-memory databases, resource allocation and data-flow computation schemes. Another example is the popular Caffeine Java caching library [19], in which a single thread is responsible for updating the internal cache data structures and meta-data. As depicted in Figure 1b, in such architectures, a single thread is responsible for each shard, in order to avoid costly synchronization while manipulating the state of a specific shard. In this case, multiple feeder threads (e.g., that communicate with external clients) insert requests into the queues according to the shards. Each thread that is responsible for a given shard then repeatedly dequeues the next request for the shard, executes it, dequeues the next request, etc. Similarly, in a data flow graph, multiple events may feed the same computational entity (e.g., a reducer that reduces the outcome of multiple mappers). Here, again, each



(a) SPMC queue used in a master worker architecture. A single master queues up tasks to be executed, which are picked by worker threads on a first comes first served basis.



(b) MPSC queues used in a sharded architecture. Here, each shard is served by a single worker thread to avoid synchronization inside the shard. Multiple collector threads can feed the queue of each shard.

Figure 1: SPMC vs. MPSC queues.

computational entity might be served by a single thread while multiple threads are sending it items, or requests, to be handled.

MPMC is the most general form of a queue and can be used in any scenario. Therefore, MPMC is also the most widely studied data structure [16, 17, 20, 24, 25, 26, 28]. Yet, this may come at a premium compared to using a more specific queue implementation.

Specifically, concurrent accesses to the same data structure require adequate concurrency control to ensure correctness. The simplest option is to lock the entire data structure on each access, but this usually dramatically reduces performance due to the sequentiality and contention it imposes [11]. A more promising approach is to reduce, or even avoid, the use of locks and replace them with *lock-free* and *wait-free* protocols that only rely on atomic operations such as *fetch-and-add* (FAA) and *compare-and-swap* (CAS), which are supported in most modern hardware architectures [10]. Wait-free implementations are particularly appealing since they ensure that each operation always terminates in a finite number of steps.

Alas, known MPMC wait-free queues suffer from large memory overheads, intricate code complexity, and low scalability. In particular, it was shown that wait-free MPMC queues require the use of a helper mechanism [4]. On the other hand, as discussed above, there are important classes of applications for which MPSC queues are adequate. Such applications could therefore benefit if a more efficient MPSC queue construction was found. This motivates studying wait-free MPSC queues, which is the topic of this paper.

Contributions In this work we present Jiffy, a fast memory efficient wait-free MPSC queue. Jiffy is unbounded in the the number of elements that can be enqueued without being dequeued (up to the memory limitations of the machine). Yet the amount of memory Jiffy consumes at any given time is proportional to the number of such items and Jiffy minimizes the use of pointers, to reduce its memory footprint.

To obtain these good properties, Jiffy stores elements in a linked list of arrays, and only allocates a new array when the last array is being filled. Also, as soon as all elements in a given array are dequeued, the array is released. This way, a typical enqueue operation requires little more than a simple FAA and setting the corresponding entry to the enqueued value and changing its status from empty to set. Hence, operations are very fast and the number of pointers is a multiple of the allocated arrays rather than the number of queued elements.

To satisfy linearizability and wait-freedom, a dequeue operation in Jiffy may return a value that is already past the head of the queue, if the enqueue operation working on the head is still on-going. To ensure correctness, we devised a novel mechanism to handle such entries both during their immediate dequeue as well as during subsequent dequeues.

Another novel idea in Jiffy is related to its buffer allocation policy. In principle, when the last buffer is full, the naive approach is for each enqueuer at that point to allocate a new buffer and then try adding it to the queue with a CAS. When multiple enqueuers try this concurrently, only one succeeds and the others need to free their allocated buffer. However, this both creates contention on the end of the queue and wastes CPU time in allocating and freeing multiple buffers each time. To alleviate these phenomena,

in Jiffy the enqueuer of the second entry in the last buffer already allocates the next buffer and tries to add it using CAS. This way, almost always, when enqueuers reach the end of a buffer, the next buffer is already available for them without any contention.

We have implemented Jiffy and evaluated its performance in comparison with three other leading lock-free and wait-free implementations, namely WFqueue [32], CCqueue [7], and MSqueue [20]. We also examined the memory requirements for the data and code of all measured implementations using valgrind [23]. The results indicate that Jiffy is up to 50% faster than WFqueue and roughly 10 times times faster than CCqueue and MSqueue. Jiffy is also more scalable than the other queue structures we tested, enabling more than 20 million operations per second even with 128 threads. Finally, the memory footprint of Jiffy is roughly 90% better than its competitors in the tested workloads, and provides similar benefits in terms of number of cache and heap accesses. Jiffy obtains better performance since the size of each queue node is much smaller and there are no auxiliary data structures. For example, in WFqueue, which also employs a linked list of arrays approach, each node maintains two pointers, there is some perthread meta-data, the basic slow-path structure (even when empty), etc. Further, WFqueue employs a lazy reclamation policy, which according to its authors is significant for its performance. Hence, arrays are kept around for some time even after they are no longer useful. In contrast, the per-node meta-data in Jiffy is just a 2-bit flag, and arrays are being freed as soon as they become empty. This translates to a more (hardware) cache friendly access pattern (as is evident in Tables 1 and 2). Also, in Jiffy dequeue operations do not invoke any atomic (e.g., FAA & CAS) operations at all.

2 Related Work

Implementing concurrent queues is a widely studied topic [2, 3, 6, 9, 18, 24, 29, 31]. Below we focus on the most relevant works.

Multi-Multi Queues: One of the most well known lock-free queue constructions was presented by Michael and Scott [20], aka *MSqueue*. It is based on a singly-linked list of nodes that hold the enqueued values plus two references to the head and the tail of the list. Their algorithm does not scale past a few threads due to contention on the queue's head and tail.

Kogan and Petrank introduced a wait-free variant of the MSqueue [14]. Their queue extends the helping technique already employed by Michael and Scott to achieve wait freedom with similar performance characteristics. They achieve wait-freedom by assigning each operation a dynamic age-based priority and making threads with younger operations help older operations to complete through the use of another data structure named a *state* array.

Morrison and Afek proposed LCRQ [22], a non-blocking queue based on a linked-list of circular ring segments, CRQ for short. LCRQ uses FAA to grab an index in the CRQ. Enqueue and dequeue operations in [22] involve a double-width compare-and-swap (CAS2).

Yang and Mellor-Crummey proposed WFqueue, a wait free queue based on FAA [32]. WFqueue utilizes a linked-list of fixed size segments. Their design employs the fast-path-slow-path methodology [15] to transform a FAA based queue into a wait-free queue. That is, an operation on the queue first tries the fast path implementation until it succeeds or the number of failures exceeds a threshold. If necessary, it falls back to the slow-path, which guarantees completion within a finite number of attempts. Each thread needs to register when the queue is started, so the number of threads cannot change after the initialization of the queue. In contrast, in our queue a thread can join anytime during the run.

Fatourou and Kallimanis proposed CCqueue [7], a blocking queue that uses combining. In CCqueue, a single thread scans a list of pending operations and applies them to the queue. Threads add their operations to the list using SWAP.

Tsigas and Zhang proposed a non-blocking queue that allows the head and tail to lag at most m nodes behind the actual head and tail of the queue [30]. Additional cyclic array queues are described in [8, 27]. Recently, a lock-free queue that extends MSqueue [20] to support batching operations was presented in [21].

Limited Concurrency Queues: David proposed a sublinear time wait-free queue [5] that supports multiple dequeuers and one enqueuer. His queue is based on infinitely large arrays. The author states that he can bound the space requirement, but only at the cost of increasing the time complexity to O(n), where n is the number of dequeuers.

Jayanti and Petrovic proposed a wait-free queue implementation supporting multiple enqueuers and

one concurrent dequeuer [13]. Their queue is based on a binary tree whose leaves are linear linked lists. Each linked list represents a "local" queue for each thread that uses the queue. Their algorithm keeps one local queue at each process and maintains a timestamp for each element to decide the order between the elements in the different queues.

3 Preliminaries

We consider a standard shared memory setting with a set of threads accessing the shared memory using only the following atomic operations:

- Store Atomically replaces the value stored in the target address with the given value.
- Load Atomically loads and returns the current value of the target address.
- CAS Atomically compares the value stored in the target address with the expected value. If those are equal, replaces the former with the desired value and the Boolean value true is returned. Otherwise, the shared memory is unchanged and false is returned.
- FAA -Atomically adds a given value to the value stored in the target address and returns the value in the target address held previously.

We assume that a program is composed of multiple such threads, which specifies the order in which each thread issues these operations and the objects on which they are invoked. In an execution of the program, each operation is invoked, known as its *invocation* event, takes some time to execute, until it terminates, known as its *termination* event. Each termination event is associated with one or more values being returned by that operation. An execution is called sequential if each operation invocation is followed immediately by the same operation's termination (with no intermediate events between them). Otherwise, the execution is said to be concurrent. When one of the operations in an execution σ is being invoked on object x, we say that x is being accessed in σ .

Given an execution σ , we say that σ induces a (partial) real-time ordering among its operations: Given two operations o_1 and o_2 in σ , we say that o_1 appears before o_2 in σ if the invocation of o_2 occurred after the termination of o_1 in σ . If neither o_1 nor o_2 appears before the other in σ , then they are considered concurrent operations in σ . Obviously, in a sequential execution the real-time ordering is a total order.

Also, we assume that each object has a sequential specification, which defines the set of allowed sequential executions on this object. A sequential execution σ is called legal w.r.t. a given object x if the restriction of this execution to operations on x (only) is included in the sequential specification of x. σ is said to be legal if it is legal w.r.t. any object being accessed in σ . Further, we say that two executions σ and σ' are equivalent if there is a one-to-one mapping between each operation in σ to an operation in σ' and the values each such pair of operations return in σ and σ' are the same.

Linearizability An execution σ is *linearizable* [12] if it is equivalent to a legal sequential execution σ' and the order of all operations in σ' respects the real-time order of operations in σ .

4 MPSC Queue

One of the problems of multi producer multi consumer queues is that any wait-free implementation of such queues requires utilizing helper data structures [4], which are both memory wasteful and slow down the rate of operations. By settling for single consumer support, we can avoid helper data structures. Further, for additional space efficiency, we seek solutions that minimize the use of pointers, as each pointer consumes 64 bits on most modern architectures. The requirement to support unbounded queues is needed since there can be periods in which the rate of enqueue operations surpasses the rate of dequeues. Without this property, during such a period some items might be dropped, or enqueuers might need to block until enough existing items are dequeued. In the latter case the queue is not wait-free.

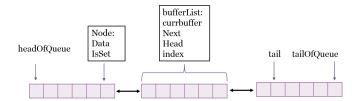


Figure 2: A node in Jiffy consists of two fields: the data itself and an IsSet field that notifies when the data is ready to be read. The BufferList consists of 5 fields: Currbuffer is an array of nodes, a Next and Prev pointers, Head index pointing to the last place in that buffer that the consumer read, and PositionInQueue that tracks the location of the BufferList in the list. The Tail index indicates the last place to be written to by the producers. HeadOfQueue is the consumer pointer for the first buffer; once the head reaches the end of the buffer, the corresponding BufferList is deleted. TailOfQueue points to the last BufferList in the linked list; once the Tail reaches the end of this buffer a new array BufferList is added.

4.1 Overview

Our queue structure is composed of a linked list of buffers. A buffer enables implementing a queue without pointers, as well as using fast FAA instructions to manipulate the head and tail indices of the queue. However, a single buffer limits the size of the queue, and is therefore useful only for bounded queue implementations. In order to support unbounded queues, we extend the single buffer idea to a linked list of buffers. This design is a tradeoff point between avoiding pointers as much as possible while supporting an unbounded queue. It also decreases the use of CAS to only adding a new buffer to the queue, which is performed rarely.

Once a buffer of items has been completely read, the (sole) consumer deletes the buffer and removes it from the linked list. Hence, deleting obsolete buffers requires no synchronization. When inserting an element to the queue, if the last buffer in the linked list has room, we perform an atomic FAA to a tail index and insert the element. Otherwise, the producer allocates a new buffer and tries to insert it to the linked list via an atomic CAS, which keeps the overall size of the queue small while supporting unbounded sizes¹. The structure of a Jiffy queue is depicted in Figure 2.

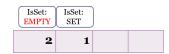
Figure 3 exhibits why the naive proposal for the queue violates linearizability. This scenario depicts two concurrent enqueue operations $enqueue_1$ and $enqueue_2$ for items i_1 and i_2 respectively, where $enqueue_1$ terminates before $enqueue_2$. Also, assume a dequeue operation that overlaps only with $enqueue_2$ (it starts after $enqueue_1$ terminates). It is possible that item i_2 is inserted at an earlier index in the queue than the index of i_1 because this is the order by which their respective FAA instructions got executed. Yet, the insertion of i_1 terminates quickly while the insertion of i_2 (into the earlier index) takes longer. Now the dequeue operation sees that the head of the queue is still empty, so it returns immediately with an empty reply. But since $enqueue_1$ terminates before the dequeue starts, this violates linearizability.



Figure 3: A linearizability problem in the basic queue idea.

Yet, the queue as described so far is not linearizable, as exemplified in Figure 3. To make the queue linearizable, if the dequeuer finds that the isSet flag is not set (nor handled, to be introduced shortly), as illustrated in Figure 4a, the dequeuer continues to scan the queue until either reaching the tail or finding an item x at position i whose isSet is set but the isSet of its previous entry is not, as depicted in Figure 4b. In the latter case, the dequeuer verifies that the isSet of all items from the head until position i are still neither set nor handled and then returns x and sets the isSet at position i to

¹In fact, to aviod contention the new buffer is usually already allocated earlier on; more accurate details appear below.





(a) The dequeuer finds that the *isSet* flag for the element pointed by the head is Empty.

(b) The dequeuer found an item whose isSet is ready, the isSet of the entry pointed by head is not set. It removes this item and marks its isSet as handled.

Figure 4: An example of the solution for the basic queue's linearizability problem.

```
Algorithm 1 Internal Jiffy classes and fields

1: class Node {
2:  T data;
3: atomic < State > isSet; }
4: // isSet has three values: empty, set and handled.
5: class bufferList {
6: Node * currBuffer;
7: atomic < bufferList* > next;
8: bufferList * prev;
9: unsigned int head;
10: unsigned int positionInQueue; }
11: bufferList * headOfQueue;
```

handled (to avoid re-dequeuing it in the future). This double reading of the *isSet* flag is done to preserve linearizability in case an earlier enqueuer has terminated by the time the newer item is found. In the latter case the dequeue should return the earlier item to ensure correct ordering. This becomes clear in the correctness proof.

4.2 Jiffy Queue Algorithm

12: atomic < bufferList* > tailOfQueue; 13: $atomic < unsigned\ int >\ tail;$

4.2.1 Internal Jiffy classes and fields

Algorithm 1 depicts the internal classes and variables. The *Node* class represents an item in the queue. It holds the data itself and a flag to indicate whether the data is ready to be read ($isSet=\mathtt{set}$), the node is empty or is in inserting process ($isSet=\mathtt{empty}$), or it was already read by the dequeuer ($isSet=\mathtt{handled}$). Every node in the queue starts with $isSet=\mathtt{empty}$.

The bufferList class represents each buffer in the queue: currbuffer is the buffer itself – an array of nodes; next is a pointer to the next buffer – it is atomic as several threads (the enqueuers) try to change it simultaneously; prev is a pointer to the previous buffer – it is not concurrently modified and therefore not atomic; head is an index to the first location in the buffer to be read by the dequeuer – it is changed only by the single dequeuer thread; positionInQueue is the position of the buffer in the queue – it is used to calculate the amount of items in the queue up until this buffer, it never changes once set and is initialized to 1.

The rest are fields of the queue: headOfQueue points to the head buffer of the queue. The dequeuer reads from the buffer pointed by headOfQueue at index head in that buffer. It is only changed by the single threaded dequeuer. tailOfQueue points to the last queue's buffer. It is atomic as it is modified by several threads. tail, initialized to 0, is the index of the last queued item. All threads do FAA on tail to receive a location of a node to insert into.

4.2.2 The Enqueue Operation

A high-level pseudo-code for enqueue operations appears in Algorithm 2; a more detailed listing appears in Appendix A. The enqueue operation can be called by multiple threads acting as producers. The method starts by performing FAA on *Tail* to acquire an index in the queue (Line 2). When there is a burst of enqueues the index that is fetched can be in a buffer that has not yet been added to the queue, or a buffer that is prior to the last one. Hence, the enqueuer needs to find in which buffer it should add its new item.

Algorithm 2 Enqueue operation

```
function ENQUEUE (data)
 3:
       while the location is in an unallocated buffer do
4:
5:
           allocate a new buffer and try to adding it to the queue with a CAS on tailOfQueue
           if unsuccessful then
6:
              delete the allocated buffer and move to the new buffer
           end if
8:
9:
       end while
       bufferList* tempTail = tailOfQueue
10:
       while the location is not in the buffer pointed by tempTail do
11:
           tempTail = tempTail \rightarrow prev
12:
        end while
13:
        //location is in this buffer
        adjust location to its corresponding index in tempTail
14:
15:
        tempTail[location].data = data
16:
        tempTail[location].isSet = set
       if location is the second entry of the last buffer then
17:
           allocate a new buffer and try adding it with a CAS on tail Of Queue; if unsuccessful, delete this buffer
18:
       end if
20: end function
                                             Head
                                                     Head
                                                                                                               Head
    EMPTY
                                              Η
                                                     EMPTY
                                                                                       EMPTY
                                                                                                      Н
                                                                                                              EMPTY
```

Figure 5: "folding" the queue - A thread fetches an index in the queue and stalls before completing the enqueue operation. Here we only keep the buffer that contains this specific index and delete the rest. H stands for $isSet = \mathtt{handled}$ and EMPTY stands for $isSet = \mathtt{empty}$.

If the index is beyond the last allocated buffer, the enqueuer allocates a new buffer and tries to add it to the queue using a CAS operation (line 3). If several threads try simultaneously, one CAS will succeed and the other enqueuers' CAS will fail. Every such enqueuer whose CAS fails must delete its buffer (line 6). If the CAS succeeds, the enqueuer moves tail Of Queue to the new buffer. If there is already a next buffer allocated then the enqueuer only moves the tail Of Queue pointer through a CAS operation.

If the index is earlier in the queue compared to the tail index (line 10), the enqueuer retracts to the previous buffer. When the thread reaches the correct buffer, it stores the data in the buffer index it fetched at the beginning (line 15), marks the location as Set (line 16) and finishes. Yet, just before returning, if the thread is in the last buffer of the queue and it obtained the second index in that buffer, the enqueuer tries to add a new buffer to the end of the queue (line 17). This is an optimization step to prevent wasteful contention prone allocations of several buffers and the deletion of most as will be explain next.

Notice that if we only had the above mechanism, then each time a buffer ends there could be a lot of contention on adding the new buffer, and many threads might allocate and then delete their unneeded buffer. This is why the enqueuer that obtains the second entry in each buffer already allocates the next one. With this optimization, usually only a single enqueuer tries to allocate such a new buffer and by the time enqueuers reach the end of the current buffer a new one is already available to them. On the other hand, we still need the ability to add a new buffer if one is not found (line 3) to preserve wait-freedom.

4.2.3 The Dequeue Operation

A high-level pseudo-code for dequeue operations appears in Algorithm 3; a more detailed listing appears in Appendix A. A dequeue operation is called by a single thread, the consumer. A dequeue starts by advancing the head index to the first element not marked with isSet = handled as such items are already dequeued (line 4). If the consumer reads an entire buffer during this stage, it deletes it (line 7). At the end of this scan, the dequeuer checks if the queue is empty and if so returns false (line 10).

Next, if the first item is in the middle of an enqueue process (isSet=empty), the consumer scans the queue (line 15) to avoid the linearizability pitfall mentioned above. If there is an element in the queue that is already set while the element pointed by the head is still empty, then the consumer needs to dequeue the latter item (denoted tempN in line 15).

Algorithm 3 Dequeue operation

```
mark the element pointed by head as n
3:
4:
5:
6:
7:
       //Skip to the first non-handled element (due to the code below, it might not be pointed by head!)
       while n.isSet == handled do
           advance n and head to the next element
           if the entire buffer has been read then
              move to the next buffer if exists and delete the previous buffer
8:
9:
       end while
10:
       if queue is empty then
11:
           return false
12:
        end if
13:
        // If the queue is not empty, but its first element is, there might be a Set element further on – find it
14:
        if n.isSet == empty then
15:
           for (tempN = n; tempN.isSet != set and not end of queue; advance tempN to next element) do
16:
              if the entire buffer is marked with handled then
                   "fold" the queue by deleting this buffer and move to the next buffer if exists
17:
18:
                                                ▷ Notice comment about a delicate garbage collection issue in the description text
19:
               end if
20:
           end for
21:
           if reached end of queue then
22:
              return false
23:
           end if
24:
       end if
25:
        // Due to concurrency, some element between n and tempN might have been set – find it
26:
        for (e = n; e = tempN \text{ or } e \text{ is before } tempN; advance } e \text{ to next element}) do
27:
           if e is not n and e.isSet == set then
28:
               tempN = e and restart the for loop again from n
29:
           end if
30:
        end for
31:
        // we scanned the path from head to tempN and did not find a prior set element - remove tempN
        tempN.isSet = \mathtt{handled}
33:
       if the entire buffer has been read then
34:
           move to the next buffer if exists and delete the previous buffer
35:
        else if tempN = n then
36:
           advance head
37:
        end if
        return tempN.data
39: end function
```

During the scan, if the consumer reads an entire buffer whose cells are all marked handled, the consumer deletes this buffer (line 17). We can think of this operation as "folding" the queue by removing buffers in the middle of the queue that have already been read by the consumer. Hence, if a thread fetches an index in the queue and stalls before completing the enqueue operation, we only keep the buffer that contains this specific index and may delete all the rest, as illustrated in Figure 5.

Further, before dequeuing, the consumer scans the path from head to tempN to look for any item that might have changed it status to **set** (line 28). If such an item is found, it becomes the new tempN and the scan is restarted.

Finally, a dequeued item is marked handled (line 32), as also depicted in Figure 4b. Also, if the dequeued item was the last non-handled in its buffer, the consumer deletes the buffer and moves the head to the next one (line 35).

There is another delicate technical detail related to deleting buffers in non-garbage collecting environments such as the run-time of C++. For clarity of presentation and since it is only relevant in non-garbage collecting environments, the following issue is not addressed in Figure 3, but is rather deferred to the detailed description in Figure 5 at the Appendix. Specifically, when preforming the fold, only the array is deleted, which consumes the most memory, while the much smaller meta-data structure of the array is transferred to a dedicated garbage collection list. The reason for not immediately deleting the the entire array's structure is to let enqueuers, who kept a pointer for tailOfQueue point to a valid bufferList at all times, with the correct prev and next pointers, until they are done with it. The exact details of this garbage collection mechanism are discussed in Appendix A.

4.2.4 Memory Buffer Pool Optimization

Instead of always allocating and releasing buffers from the operating system, we can maintain a buffer pool. This way, when trying to allocate a buffer, we first check if there is already a buffer available in the buffer pool. If so, we simply claim it without invoking an OS system call. Similarly, when releasing a buffer, rather than freeing it with an OS system call, we can insert it into the buffer pool. It is possible

to have a single shared buffer pool for all threads, or let each thread maintain its own buffer pool. This optimization can potentially reduce execution time at the expense of a somewhat larger memory heap area. The code used for Jiffy's performance measurements does *not* employ this optimization. In Jiffy, allocation and freeing of buffers are in any case a relatively rare event.

5 Correctness

5.1 Linearizability

To prove linearizability, we need to show that for each execution σ that may be generated by Jiffy, we can find an equivalent legal sequential execution that obeys the real-time order of operations in σ . We show this by constructing such an execution. Further, for simplicity of presentation, we assume here that each value can be enqueued only once. Consequently, it is easy to verify from the code that each value can also be returned by a dequeue operation at most once. In summary we have:

Observation 5.1. Each enquequed value can be returned by a dequeue operation at most once.

Theorem 5.2. The Jiffy queue implementation is linearizable.

Proof. Let σ be an arbitrary execution generated by Jiffy. We now build an equivalent legal sequential execution σ' . We start with an empty sequence of operations σ' and gradually add to it all operations in σ until it is legal and equivalent to σ . Each operation inserted into σ' is collapsed such that its invocation and termination appear next to each other with no events of other operations in between them, thereby constructing σ' to be sequential.

First, all dequeue operations of σ are placed in σ' in the order they appear in σ . Since Jiffy only supports a single dequeuer, then in σ there is already a total order on all dequeue operations, which is preserved in σ' .

Next, by the code, a dequeue operation deq that does not return empty, can only return a value that was inserted by an enqueue operation enq that is concurrent or prior to deq (only enqueue operations can change an entry to set). By Observation 5.1, for each such deq operation there is exactly one such enq operation. Hence, any ordering in which enq appears before deq would preserve the real-time ordering between these two operations.

Denote the set of all enqueue operations in σ by ENQ and let \widehat{ENQ} be the subset of ENQ consisting of all enqueue operations enq such that the value enqueued by enq is returned by some operation deq in σ . Next, we order all enqueue operations in \widehat{ENQ} in the order of the dequeue operations that returned their respective value.

Claim 5.3. The real time order in σ is preserved among all operations in \widehat{ENQ} .

Proof of Claim 5.3. Suppose Claim 5.3 does not hold. Then there must be two enqueue operations enq_1 and enq_2 and corresponding dequeue operations deq_1 and deq_2 such that the termination of enq_1 is before the invocation of enq_2 in σ , but deq_2 occurred before deq_1 . Denote the entry accessed by enq_1 by in_1 and the entry accessed by enq_2 by in_2 . Since the invocation of enq_2 is after the termination of enq_1 , then during the invocation of enq_2 the status of in_1 was already set (from line 16 in Algorithm 2). Moreover, in_2 is further in the queue than in_1 (from line 2 in Algorithm 2). Since deq_2 returned the value enqueued by enq_2 , by the time deq_2 accessed in_2 the status of in_1 was already set. As we assumed deq_1 is after deq_2 , no other dequeue operation has dequeued the value in in_1 . Hence, while accessing in_2 and before terminating, deq_2 would have checked in_1 and would have found that it is now set (the loops in lines 4 and 26 of Algorithm 3 – this is the reason for line 26) and would have returned that value instead of the value at in_2 . A contradiction.

To place the enqueue operations of \widehat{ENQ} inside σ' , we scan them in their order in \widehat{ENQ} (as defined above) from earliest to latest. For each such operation enq that has not been placed yet in σ : (i) let deq' be the latest dequeue operation in σ that is concurrent or prior to enq in σ and neither deq' nor any prior dequeue operation return the value enqueued by enq, and (ii) let deq'' be the following dequeue operation in σ ; we add enq in the last place just before deq'' in σ' . This repeats until we are done placing all operations from \widehat{ENQ} into σ' .

Claim 5.4. The real-time ordering in σ between dequeue operations and enqueue operations in \widehat{ENQ} is preserved in σ' as built thus far.

Proof of Claim 5.4. Since we placed each enqueue operation enq before any dequeue operation whose invocation is after the termination of enq, we preserve real-time order between any pair of enqueue and dequeue operations. Similarly, by construction the relative order of dequeue operations is not modified by the insertion of enqueue operations into σ' . Thus, real-time ordering is preserved.

Hence, any potential violation of real time ordering can only occur by placing enqueue operations in a different order (with respect to themselves) than they originally appeared in \widehat{ENQ} . For this to happen, it means that there are two enqueue operation enq_1 and enq_2 such that enq_1 appears before enq_2 in \widehat{ENQ} but ended up in the reverse order in σ' . Denote deq'_1 the latest dequeue operation in σ that is prior or concurrent to enq_1 and neither deq'_1 nor any prior dequeue operation return the value enqueued by enq_1 and similarly denote deq'_2 for enq_2 . Hence, if enq_2 was inserted at an earlier location than enq_1 , then deq'_2 is also before deq'_1 . But since the order of enq_1 and enq_2 in \widehat{ENQ} preserves their real time order, the above can only happen if the value enqueued by enq_2 was returned by an earlier dequeue than the one returning the value enqueued by enq_1 . Yet, this violates the definition of the ordering used to create \widehat{ENQ} .

Claim 5.5. The constructed execution σ' preserves legality.

Proof of Claim 5.5. By construction, enqueue operations are inserted in the order their values have been dequeued, and each enqueue is inserted before the dequeue that returned its value. Hence, the only thing left to show is legality w.r.t. dequeue operations that returned empty.

To that end, given a dequeue operation deq_i that returns empty, denote $\#deq_{\sigma',i}$ the number of dequeue operations that did not return empty since the last previous dequeue operation that did return empty, or the beginning of σ' of none exists. Similarly, denote $\#enq_{\sigma',i}$ the number of enqueue operations during the same interval of σ' .

Sub-Claim 5.6. For each deq_i that returns empty, $\#\widehat{deq}_{\sigma',i} > \#enq_{\sigma',i}$.

Proof of Sub-Claim 5.6. Recall that the ordering in σ' preserves the real time order w.r.t. σ and there is a single dequeuer. Assume by way of contradiction that the claim does not hold, and let deq_i be the first dequeue operation that returned empty while $\#\widehat{deq}_{\sigma',i} \leq \#enq_{\sigma',i}$. Hence, there is at least one enqueue operation enq_j in the corresponding interval of σ' whose value is not dequeued in this interval. In this case, deq_i cannot be concurrent to enq_j in σ since otherwise by construction enq_j would have been placed after it (as it does not return its value). Hence, deq_i is after enq_j in σ . Yet, since each dequeue removes at most one item from the queue, when deq_i starts, the tail of the queue is behind the head and there is at least one item whose state is set between them. Thus, by lines 4 and 26 of Algorithm 3, deq_i would have returned one of these items rather than return empty. A contradiction.

With Sub-Claim 5.6 we conclude the proof that σ' as constructed so far is legal.

The last thing we need to do is to insert enqueue operations whose value was not dequeued, i.e., all operations in $\widehat{ENQ} \setminus \widehat{ENQ}$. Denote by enq' the last operation in \widehat{ENQ} .

Claim 5.7. Any operation $enq'' \in ENQ \setminus \widehat{ENQ}$ is either concurrent to or after enq' in σ .

Proof of Claim 5.7. Assume, by way of contradiction, that there is an operation $enq'' \in ENQ \setminus \widehat{ENQ}$ that is before enq' in σ . Hence, by the time enq' starts, the corresponding entry of enq'' is already set (line 16 of Algorithm 2) and enq' obtains a later entry (line 2). Yet, since dequeue operations scan the queue from head to tail until finding a set entry (lines 4, 15, and 26 in Algorithm 3), the value enqueued by enq'' would have been dequeued before the value of enq'. A contradiction.

Hence, following Claim 5.7, we insert to σ' all operations in $ENQ \setminus \widehat{ENQ}$ after all operations of \widehat{ENQ} . In case of an enqueue operation enq that is either concurrent with or later than a dequeue deq in σ , then enq is inserted to σ' after deq. Among concurrent enqueue operations, we break symmetry arbitrarily. Hence, real-time order is preserved in σ' .

	Jiffy	WF		LCRQ		CC		MS	
	Absolute	Absolute	Relative	Absolute	Relative	Absolute	Relative	Absolute	Relative
Total Heap Usage	38.70 MB	611.63 MB	x15.80	610.95 MB	x15.78	305.18 MB	x7.88	1.192 GB	x31.54
Number of Allocs	3,095	9,793	x3.16	1,230	x0.40	5,000,015	x1,615	5,000,010	x1,615
Peak Heap Size	44.81 MB	200.8 MB	x4.48	612.6 MB	x13.67	420.0 MB	x9.37	1.229 GB	x28.08
# of Instructions Executed	550,416,453	5,612,941,764	x10.20	1,630,746,827	x2.96	3,500,543,753	x6.36	1,821,777,428	x3.31
I1 Misses	2,162	1,714	x0.79	1,601	x0.74	1,577	x0.73	1,636	x0.76
L3i Misses	2,084	1,707	x0.82	1,591	x0.76	1,572	x0.75	1,630	x0.78
Data Cache Tries (R+W)	281,852,749	2,075,332,377	x7.36	650,257,906	x2.3	1,238,761,631	x4.40	646,304,575	x2.29
D1 Misses	1,320,401	25,586,262	x19.37	20,037,956	x15.17	15,000,507	x11.36	11,605,064	x8.79
L3d Misses	652,194	10,148,090	x15.56	5,028,204	x7.7	14,971,182	x22.96	10,973,055	x16.82

Table 1: Valgrind memory usage statistics run with one enqueuer and one dequeuer. I1/D1 is the L1 instruction/data cache respectively while L3i/L3d is the L3 instruction/data cache respectively.

The only thing to worry about is the legality of dequeue operations that returned empty. For this, we can apply the same arguments as in Claim 5.5 and Sub-claim 5.6. In summary, σ' is an equivalent legal sequential execution to σ that preserves σ 's real-time order.

5.2 Wait-Freedom

We show that each invocation of enqueue and dequeue returns in a finite number of steps.

Lemma 5.8. Each enqueue operation in Algorithm 2 completes in a finite number of steps.

Proof. An enqueue operation, as listed in Algorithm 2, consists of two while loops (line 3 and line 10), each involving a finite number of operations, a single FAA operation at the beginning (line 2), and then a short finite sequence of operations from line 15 onward. Hence, we only need to show that the two while loops terminate in a finite number of iterations.

The goal of the first while loop is to ensure that the enqueuer only accesses an allocated buffer. That is, in each iteration, if the *location* index it obtained in line 2 is beyond the last allocated buffer, the enqueuer allocates a new buffer and tries to add it with a CAS to the end of the queue (line 3). Even if there are concurrent such attempts by multiple enqueuers, in each iteration at least one of them succeeds, so this loop terminates in a finite number of steps at each such enqueuer.

The next step for the enqueuer is to obtain the buffer corresponding to *location*. As mentioned before, it is possible that by this time the queue has grown due to concurrent faster enqueue operations. Hence, the enqueuer starts scanning from tailOfQueue backwards until reaching the correct buffer (line 10). Since new buffers are only added at the end of the queue, this while loop also terminates in a finite number of steps regardless of concurrency.

Lemma 5.9. Each dequeue operation in Algorithm 3 completes in a finite number of steps.

Proof. We prove the lemma by analyzing Algorithm 3. Consider the while loop at line 4. Here, we advance head to point to the first non-handled element, in case it is not already pointing to one. As indicated before, the latter could occur if a previous dequeue deq_1 removed an element not pointed by head. As we perform this scan only once and the queue is finite, the loop terminates within a finite number of iterations, each consisting of at most a constant number of operations.

If at this point the queue is identified as empty, which is detectable by comparing head and tail pointers and indices, then the operation returns immediately (line 10). The following for loop is at line 15. Again, since the queue is finite, the next set element is within a finite number of entries away. Hence, we terminate the for loop after a finite number of iterations.

The last for loop iteration at line 26 scans the queue from its head to the first set element, which as mentioned before, is a finite number of entries away. Yet, the for loop could be restarted in line 28, so we need to show that such a restart can occur only a finite number of times. This is true because each time we restart, we shorten the "distance" that the for loop at line 26 needs to cover, and as just mentioned, this "distance" is finite to begin with.

The rest of the code from line 32 onward is a short list of simple instructions. Hence, each dequeue operation terminates in a finite number of steps. \Box

Theorem 5.10. The Jiffy queue implementation is wait-free.

6 Performance Evaluation

We compare Jiffy to several representative queue implementations in the literature: Yang and Mellor-Crummey's queue [32] is the most recent wait free FAA-based queue, denoted WFqueue; Morrison and Afek LCRQ [22] as a representative of nonblocking FAA-based queues; Fatourou and Kallimanis CCqueue [7] is a blocking queue based on the combining principle. We also test Michael and Scott's classic lock-free MSqueue [20]. We include a microbenchmark that only preforms FAA on a shared variable. This serves as a practical upper bound for the throughput of all FAA based queue implementations. Notice that Jiffy performs FAA only during enqueue operations, but not during dequeues. Also, measurements of Jiffy do not include a memory buffer pool optimization mentioned in Section 4.2.4.

Implementation: We implemented our queue algorithm in C++ [1]. We compiled Jiffy with g++ version 7.4.0 with -Os optimization level. We use the C implementation provided by [32] for the rest of the queues mentioned here. They are compiled with GCC 4.9.2 with -Os optimization level. The buffer size of our queue is 1620 entries. The segment size of Yang and Mellor-Crummey queue is 2^{10} and in LCRQ it is 2^{12} , the optimal sizes according to their respective authors.

Platforms: We measured performance on the following servers:

- AMD PowerEdge R7425 server with two AMD EPYC 7551 Processors. Each processor has 32 2.00GHz/2.55GHz cores, each of which multiplexes 2 hardware threads, so in total this system supports 128 hardware threads. The hardware caches include 32K L1 cache, 512K L2 cache and 8192K L3 cache, and there are 8 NUMA nodes, 4 per processor.
- Intel Xeon E5-2667 v4 Processor including 8 3.20GHz cores with 2 hardware threads, so this system supports 16 hardware threads. The hardware caches include 32K L1 cache, 256K L2 cache and 25600K L3 cache.

Methodology: We used two benchmarks: one that only inserts elements to the queue (enqueue only benchmark) whereas the second had one thread that only dequeued while the others only enqueued.

In each experiment, x threads concurrently perform operations for a fixed amount of seconds, as specified shortly. We tightly synchronize the start and end time of threads by having them spin-wait on a "start" flag. Once all threads are created, we turn this flag on and start the time measurement. To synchronize the end of the tests, each thread checks on every operation an "end" flag (on a while loop). When the time we measure ends we turn this flag on. Each thread then counts the amount of finished operations it preformed and all are combined with FAA after the "end" flag is turned on.

In order to understand the sensitivity of our results to the run lengths, we measure both 1 and 10 seconds runs. That is, the fixed amount of time between turning on the "start" and "end" flags is set to 1 and 10 seconds, respectively. The graphs depict the throughput in each case, i.e., the number of operations applied to the shared queue per second by all threads, measured in million operations per second (MOPS). Each test was repeated 11 times and the average throughput is reported. All experiments employ an initially empty queue.

With AMD PowerEdge for all queues we pinned each software thread to a different hardware thread based on the architecture of our server. We started mapping to the first NUMA node until it became full, with two threads running on the same core. We continue adding threads to the closest NUMA node until we fill all the processor hardware threads. Then we move to fill the next processor in the same way.

For Intel Xeon server we test up to 32 threads without pinning whereas the server only has 16 hardware threads. This results in multiple threads per hardware thread.

Total Space Usage We collected memory usage statistics via valgrind version 3.13.0 [23]. We measure the memory usage when inserting 10^7 elements to the queue on AMD PowerEdge. Table 1 lists the memory usage when the queues are used by one enqueuer and one dequeuer. Jiffy's memory usage is significantly smaller than all other queues compared in this work. Jiffy uses 38.7 MB of heap memory which is 93.67% less than the WFqueue consumption of 611.63 MB, and 97% less than MSqueue. The peak heap size depicts the point where memory consumption was greatest. As can be seen Jiffy's peak

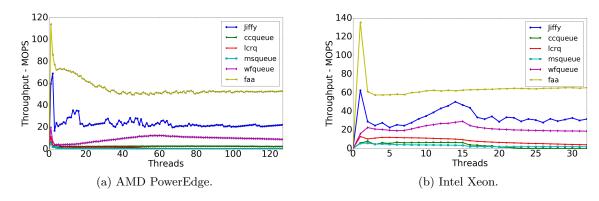


Figure 6: Enqueues only - 1 second runs.

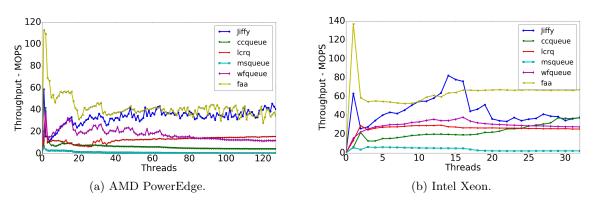


Figure 7: Multiple enqueuers with a single dequeuer - 1 second runs.

is significantly lower than the rest. Jiffy's miss ratios in L1 and L3 data caches are better due to its construction as a linked list of arrays.

Memory Usage with 128 Threads Table 2 lists the Valgrind statistics for the run with 127 enqueuers and one dequeuer. Here, Jiffy's heap consumption has grown to 77.10 MB, yet it is 87% less than WFqueue's consumption, 87% less than CCqueue and 96.8% less than MSqueue. As mentioned above, this memory frugality is an important factor in Jiffy's cache friendliness, as can be seen by the cache miss statistics of the CPU data caches (D1 and L3d miss).

	Jiffy	WF		LCRQ		CC		MS	
	Absolute	Absolute	Relative	Absolute	Relative	Absolute	Relative	Absolute	Relative
Total Heap Usage	77.10 MB	611.70 MB	x7.93	1.19 GB	x15.8	605.68 MB	x7.85	2.36 GB	x31.42
Number of Allocs	6409	10045	x1.57	2819	x0.44	9922394	x1548	9922263	x1548
Peak Heap Size	87.42 MB	624.5 MB	x7.14	1.191 GB	x13.94	796.7 MB	x9.11	2.426 GB	x28.42
# of Instructions Executed	678,651,794	3,099,458,758	x4.57	1,671,748,656	x0.99	8,909,842,602	x13.13	11,944,994,600	x17.60
I1 Misses	2,238	1,724	x0.77	1,669	x0.75	1,668	x0.75	1,689	x0.75
L3i Misses	2,194	1,717	x0.78	1,658	x0.76	1,664	x0.76	1,684	x0.77
Data Cache Tries (R+W)	352,980,193	1,849,850,595	x5.24	690,008,884	x1.95	3,172,272,116	x8.99	3,237,610,091	x9.17
D1 Misses	2,643,266	51,230,800	x19.38	30,012,317	x11.35	68,349,604	x25.86	79,097,745	x29.92
L3d Misses	1,298,646	40,453,921	x31.15	19,946,542	x15.36	57,287,592	x44.11	64,219,733	x49.45

Table 2: Valgrind memory usage statistics with 127 enqueuers and one dequeuer. I1/D1 is the L1 instruction/data cache respectively while L3i/L3d is the L3 instruction/data cache respectively.

Throughput Results Figure 6a shows results for the enqueues only benchmark on the AMD PowerEdge server for one second runs. Jiffy obtains the highest throughput with two threads, as the two threads are running on the same core and different hardware threads. The rest of the queues obtain the highest throughput with only one thread. Jiffy outperforms all queues, reaching as much as 69 millions operations per second (MOPS) with two threads. The FAA is an upper-bound for all the queues as they all preform FAA in each enqueue operation. The peak at 16 threads and the drop in 17 threads in Jiffy

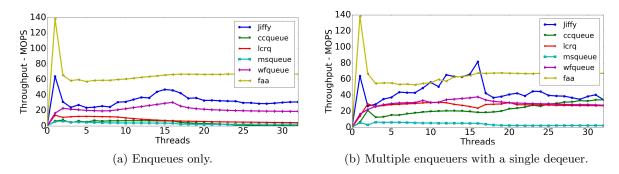


Figure 8: Intel Xeon - 10 seconds runs.

is due to the transition into a new NUMA node. At this point, a single thread is running by itself on a NUMA node. Notice also the minor peaks when adding a new core vs. starting the 2nd thread on the same core.

Beyond 64 threads the application already spans two CPUs. The performance reaches steady-state. This is because with 2 CPUs, the sharing is at the memory level, which imposes a non-negligible overhead. Jiffy maintains its ballpark performance even when 128 enqueuers are running with a throughput of 22 MOPS

Figure 6b shows results for the enqueues only benchmark on the Intel Xeon server for one second runs. All of the queues suffer when there are more threads than hardware threads, which happens beyond 16 threads.

Figure 7a shows results with a single dequeuer and multiple enqueuers on the AMD PowerEdge server for one second runs. Jiffy is the only queue whose throughput improves in the entire range of the graph (after the initial 2-threads drop).

Figure 7b shows results with a single dequeuer and multiple enqueuers on the Intel Xeon server for one second runs. Here Jiffy outperforms in some points even the FAA benchmark. This is because the dequeuer of Jiffy does not preform any synchronization operations. None of the other queues can achieve this due to their need to support multiple dequeuers.

Figure 8a shows results for the enqueues only benchmark on the Intel Xeon server when the run length is set to 10 seconds. Figure 8b shows results for a single dequeuer and multiple enqueuers on the Intel Xeon server with 10 seconds runs. As can be seen, in both cases the results are very similar to the 1 second runs. The same holds for the AMD PowerEdge server.

Let us comment that in a production system, the enqueque rate cannot surpass the dequeue rate for long periods of time; otherwise the respective queue would grow arbitrarily. However, it is likely to have short bursts lasting a few seconds each, where a single shard gets a disproportionate number of enqueues. This test exemplifies Jiffy's superior ability to overcome such periods of imbalance.

7 Conclusions

In this paper we presented Jiffy, a fast memory efficient wait-free multi-producers single-consumer FIFO queue. Jiffy is based on maintaining a linked list of buffers, which enables it to be both memory frugal and unbounded. Most enqueue and dequeue invocations in Jiffy complete by performing only a few atomic operations.

Further, the buffer allocation scheme of Jiffy is designed to reduce contention and memory bloat. Reducing the memory footprint of the queue means better inclusion in hardware caches and reduced resources impact on applications.

Jiffy outperforms prior queues in all concurrency levels especially when the dequeuer is present. Moreover, Jiffy's measured memory usage is significantly smaller than all other queues tested in this work, $\approx 90\%$ lower than WFqueue, LCRQ, CCqueue, and MSqueue.

References

[1] D. Adas. Jiffy's C++ Implementation. https://github.com/DolevAdas/Jiffy, 2020.

- [2] D. Alistarh, J. Kopinsky, J. Li, and N. Shavit. The Spraylist: A Scalable Relaxed Priority Queue. In Proc. of the ACM PPoPP, pages 11–20, 2015.
- J. Aspnes and M. Herlihy. Wait-Free Data Structures in the Asynchronous PRAM Model. In Proc. of ACM SPAA, pages 340–349, 1990.
- [4] K. Censor-Hillel, E. Petrank, and S. Timnat. Help! In Proc. of ACM PODC, pages 241–250, 2015.
- [5] M. David. A Single-Enqueuer Wait-Free Queue Implementation. In International Symposium on Distributed Computing, pages 132–143. Springer, 2004.
- [6] T. David, A. Dragojevic, R. Guerraoui, and I. Zablotchi. Log-Free Concurrent Data Structures. In USENIX ATC, pages 373–386, 2018.
- [7] P. Fatourou and N. D. Kallimanis. Revisiting the Combining Synchronization Technique. In Proc. of the ACM PPoPP, pages 257–266, 2012.
- [8] J. Giacomoni, T. Moseley, and M. Vachharajani. FastForward for Efficient Pipeline Parallelism: a Cache-Optimized Concurrent Lock-Free Queue. In *Proc. of the ACM PPoPP*, pages 43–52, 2008.
- [9] M. Hedayati, K. Shen, M. L. Scott, and M. Marty. Multi-Queue Fair Queuing. In USENIX Annual Technical Conference (ATC), pages 301–314, 2019.
- [10] M. Herlihy. Wait-Free Synchronization. ACM Transactions on Programming Languages and Systems (TOPLAS), 13(1):124–149, 1991.
- [11] M. Herlihy and N. Shavit. The Art of Multiprocessor Programming. Morgan Kaufmann, 2011.
- [12] M. P. Herlihy and J. M. Wing. Linearizability: A Correctness Condition for Concurrent Objects. ACM Transactions on Programming Languages and Systems (TOPLAS), 12(3):463–492, 1990.
- [13] P. Jayanti and S. Petrovic. Logarithmic-Time Single Deleter, Multiple Inserter Wait-Free Queues and Stacks. In *International Conference on Foundations of Software Technology and Theoretical Computer Science*, pages 408–419. Springer, 2005.
- [14] A. Kogan and E. Petrank. Wait-Free Queues with Multiple Enqueuers and Dequeuers. In Proc. of the ACM PPoPP, pages 223–234, 2011.
- [15] A. Kogan and E. Petrank. A Methodology for Creating Fast Wait-Free Data Structures. In Proc. of the ACM PPoPP, pages 141–150, 2012.
- [16] E. Ladan-Mozes and N. Shavit. An Optimistic Approach to Lock-Free FIFO Queues. In Distributed Computing, 18th International Conference (DISC), pages 117–131, 2004.
- [17] N. M. Lê, A. Guatto, A. Cohen, and A. Pop. Correct and Efficient Bounded FIFO Queues. In 25th International Symposium on Computer Architecture and High Performance Computing, pages 144–151. IEEE, 2013.
- [18] X. Liu and W. He. Active Queue Management Design Using Discrete-Event Control. In 46th IEEE Conference on Decision and Control, pages 3806–3811, 2007.
- [19] B. Manes. Caffeine: A High Performance Caching Library for Java 8. https://github.com/ben-manes/caffeine, 2017.
- [20] M. M. Michael and M. L. Scott. Simple, Fast, and Practical Non-Blocking and Blocking Concurrent Queue Algorithms. In ACM PODC, pages 267–275, 1996.
- [21] G. Milman, A. Kogan, Y. Lev, V. Luchangco, and E. Petrank. BQ: A Lock-Free Queue with Batching. In Proceedings of the ACM SPAA, pages 99–109, 2018.
- [22] A. Morrison and Y. Afek. Fast Concurrent Queues for x86 Processors. In Proc. of the ACM PPoPP, pages 103–112, 2013.
- [23] N. Nethercote and J. Seward. Valgrind: A Program Supervision Framework. *Electronic notes in theoretical computer science*, 89(2):44–66, 2003.
- [24] W. N. Scherer, D. Lea, and M. L. Scott. Scalable Synchronous Queues. In Proceedings of the ACM PPOPP, pages 147–156, 2006.
- [25] M. L. Scott. Non-Blocking Timeout in Scalable Queue-Based Spin Locks. In Proceedings of ACM PODC, pages 31–40, 2002.
- [26] M. L. Scott and W. N. Scherer. Scalable Queue-Based Spin Locks with Timeout. In Proceedings of the ACM PPOPP, pages 44–52, 2001.
- [27] N. Shafiei. Non-Blocking Array-Based Algorithms for Stacks and Queues. In Proceedings of ICDCN, pages 55–66. Springer, 2009.
- [28] N. Shavit and A. Zemach. Scalable Concurrent Priority Queue Algorithms. In Proceedings of the ACM PODC, pages 113–122, 1999.
- [29] F. Strati, C. Giannoula, D. Siakavaras, G. I. Goumas, and N. Koziris. An Adaptive Concurrent Priority Queue for NUMA Architectures. In Proceedings of the 16th ACM International Conference on Computing Frontiers (CF), pages 135–144, 2019.
- [30] P. Tsigas and Y. Zhang. A Simple, Fast and Scalable Non-Blocking Concurrent FIFO Queue for Shared Memory Multiprocessor Systems. In ACM SPAA, pages 134–143, 2001.
- [31] H. Vandierendonck, K. Chronaki, and D. S. Nikolopoulos. Deterministic Scale-Free Pipeline Parallelism with Hyperqueues. In Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, page 32. ACM, 2013.
- [32] C. Yang and J. Mellor-Crummey. A Wait-Free Queue as Fast as Fetch-and-Add. Proc. of the ACM PPoPP, 51(8):16, 2016.

A Detailed Pseudocode for Jiffy

Enqueue Operation Algorithm 4 lists the detailed implementation of the enqueue operation, expanding the high level pseudo-code given in Algorithm 2. Specifically, if the location index obtained through the fetch_add(1) operation in line 2 is beyond the last allocated buffer, the enqueuer allocates a new buffer and tries to add it to the queue using a CAS operation. This is performed by the while loop at line 6 of Algorithm 4, which corresponds to the while loop in line 3 of Algorithm 2. If the CAS fails, it means that another thread succeeded, so we can move to the new allocated buffer and check if the location we fetched at line 2 is in that buffer. If it is not, then we continue trying to add buffers until reaching the correct buffer.

The while loop at line 22 of Algorithm 4 matches the loop at line 10 in Algorithm 2, where if needed, the enqueuer retracts from the end of the queue to the previous buffer corresponding to location. Line 29 in Algorithm 4 is where we calculate the correct index. To do that we remove from location, which is a global index for the entire queue, the amount of items up to the current buffer. This is done at line 21 of Algorithm 4. We then insert the node in line 31. Yet, just before returning, if the thread is in the last buffer of the queue and it obtained the second index in that buffer, then the enqueuer tries to add a new buffer to the end of the queue at line 33.

Algorithm 4 Enqueue operation - detailed

```
1: function ENQUEUE(data)
        unsigned int location = tail.fetch\_add(1)
 3:
        bool isLastBuffer = true
bufferList* tempTail = tailOfQueue.load()
 4:
 5:
6:
        unsigned int numElements = bufferSize*tempTail → positionInQueue
        while location ≥ numElements do
             //location is in the next buffer
 8:
9:
            \mathbf{if} (tempTail \rightarrow next).load() == NULL then
                 //buffer not yet exist in the queue
10:
                 bufferList* newArr = new bufferList(bufferSize, tempTail\rightarrow positionInQueue + 1, tempTail)
11:
                 if CAS(&(tempTail→ next), NULL, newArr) then
                     CAS(&tailOfQueue,&tempTail, newArr)
13:
14:
                    {\it delete}\ {\it newArr}
15:
                 end if
16:
             end if
             tempTail = tailOfQueue.load()
17:
18:
             numElements = bufferSize*tempTail \rightarrow positionInQueue
19:
20:
21:
         //calculating the amount of item in the queue - the current buffer
         unsigned int prevSize= bufferSize*(tempTail → positionInQueue-1)
22:
         \mathbf{while} \ \mathrm{location} \ \mathsf{;} \ \mathrm{prevSize} \ \mathbf{do}
23:
             // location is in a previous buffer from the buffer pointed by tail
24:
             tempTail = tempTail→ prev
25:
             prevSize = bufferSize*(tempTail → positionInQueue - 1)
26:
             isLastBuffer = false
27:
28:
         end while
          / location is in this buffer
29:
         \overset{.}{N}ode^* n = \&(tempTail \rightarrow currbuffer[location - prevSize])
30:
        if n \rightarrow isSet.load() == State.empty then
            n \rightarrow data = data
32:
                \rightarrow isSet.store(State.set)
33:
             if index == 1 && isLastBuffer then
34:
35:
                 //allocating a new buffer and adding it to the queue bufferList* newArr = new bufferList(bufferSize, tempTail) \rightarrow positionInQueue + 1, tempTail)
36:
                 if !CAS(&(tempTail→ next), NULL, newArr) then
                    delete newArr
38:
39:
             end if
40:
         end if
41: end function
```

Dequeue Operation Algorithm 5 provides the detailed implementation for the dequeue operation, corresponding to the high level pseudo-code in Algorithm 3. First, we skip handled elements in lines 3–10 of Algorithm 5, which match the loop at line 4 of Algorithm 3. Next, we check whether the queue is empty and if so return false (lines 12–14 of Algorithm 5). To do so, we compare the headOfQueue and tailOfQueue to check if they point to the same BufferList as well as compare the tail and the head. Note that the tail index is global to the queue while the head index is a member of the BufferList class.

If the element pointed by head is marked set, we simply remove it from the queue and return (lines 15–20). Next, if the first item is in the middle of an enqueue process (isSet=empty), then the consumer scans the queue for a later set item. This is performed by function SCAN listed in Algorithm 8, which matches line 15 in Algorithm 3. If such an item is found, it is marked tempN. The folding of the queue invoked in line 17 of Algorithm 3 is encapsulated in the FOLD function listed in Algorithm 6.

When preforming the fold, only the array is deleted, which consumes the most memory. The reason for not deleting the the entire array's structure is to let enqueuers, who kept a pointer for tailOfQueue at line 4 of Algorithm 4, point to a valid BufferList at all times with the correct prev and next pointers. To delete the rest of the array structure later on, the dequeuer keeps this buffer in a list called garbageList, a member of the Jiffy class (line 54). When a BufferList from the queue is deleted, the dequeuer checks if there is a BufferList in garbageList that is before the buffer being deleted. If so, the dequeuer deletes it as well (lines 70- 75).

Before dequeuing, the consumer scans the path from head to tempN to look for any item that might have changed its status to set. This is executed by function Rescan in Algorithm 9, corresponding to line 28 of Algorithm 3. If such an item is found, it becomes the new tempN and the scan is restarted. This is to check whether there is an even closer item to n that changed its status to set. Finally, a dequeued item is marked handled in line 33 of Algorithm 5, matching line 32 of Algorithm 3. Also, if the dequeued item was the last non-handled in its buffer, the consumer deletes the buffer and moves the head to the next one. This is performed in line 36 of Algorithm 5, corresponding to line 35 of Algorithm 3.

Algorithm 5 Dequeue operation - detailed

```
1: function DEQUEUE(T&data)
          Node* n = \&(headOfQueue \rightarrow currbuffer[headOfQueue \rightarrow head]);
 2:
3:
          while n \rightarrow isSet.load() == State.handled do // find first non-handled item
 4:
5:
6:
7:
              headOfQueue \rightarrow head++
              bool res = MoveToNextBuffer
() // if at end of buffer, skip to next one
              if !res then
                   return false; // reached end of queue and it is empty
 8:
9:
              end if
              n = \&(headOfQueue \rightarrow currbuffer[headOfQueue \rightarrow head]) \ // \ n \ points \ to \ the \ beginning \ of \ the \ queue
10:
11:
          end while
             / check if the queue is empty
          // check if the queue is empty if ((headOfQueue \rightarrow head = tail.load() % bufferSize )) then
12:
13:
               return false
14:
15:
          if n \to isSet.load() == State.set then // if the first element is set, dequeue and return it
               headOfQueue→head++
MoveToNextBuffer()
16:
17:
18:
               data = n→data
19:
               return true
20:
           end if
21:
22:
23:
          if n \rightarrow isSet.load() == State.empty then // otherwise, scan and search for a set element bufferList* tempHeadOfQueue = headOfQueue
               unsigned int tempHead = headOfQueue\rightarrowhead
24:
25:
               \label{eq:node_node} $\operatorname{Node}^* \operatorname{tempN} = \&(\operatorname{tempHeadOfQueue} \to \operatorname{currbuffer[tempHead]})$ bool res = Scan(\operatorname{tempHeadOfQueue}, \operatorname{tempHead}, \operatorname{tempN})
26:
               if !res then
27:
                   return false; // if none was found, we return empty
28:
29:
30:
               //here tempN == set (if we reached the end of the queue we already returned false) Rescan(headOfQueue ,tempHeadOfQueue ,tempHead ,tempN) // perform the rescan
31:
               // tempN now points to the first set element – remove tempN data = tempN \rightarrow data
32:
33:
               tempN \ris Set.store(State.handled)
34:
               if (tempHeadOfQueue==headOfQueue && tempHead ==head) //tempN ==n then
                    \begin{array}{l} {\rm headOfQueue} {\rightarrow} {\rm head} {+} {+} \\ {\rm MoveToNextBuffer}() \end{array}
35:
36:
37:
               end if
38:
               return true
          end if
40: end function
```

Algorithm 6 Folding a fully handled buffer in the middle of the queue

```
41: function FOLD(bufferList* tempHeadOfQueue,
                                                                unsigned int& tempHead, bool& flag_moveToNewBuffer,
    flag_bufferAllHandeld)
42:
        if tempHeadOfQueue == tailOfQueue.load() then
43:
            return false // the queue is empty – we reached the tail of the queue
44:
45:
         bufferList* next = tempHeadOfQueue \rightarrow next.load()
        bufferList* prev = tempHeadOfQueue→prev
if next == NULL then
46:
47:
48:
            return false // we do not have where to move
49:
        end if
50:
         // shortcut this buffer and delete it
51:
        next \rightarrow prev = prev
52:
53:
        _{\mathrm{prev} \rightarrow \mathrm{next.store}(\mathrm{next})}
        {\tt delete[] \ tempHeadOfQueue} {\to} {\tt currbuffer}
        garbageList.addLast(tempHeadOfQueue)
tempHeadOfQueue = next
54:
55:
         tempHead = tempHeadOfQueue \rightarrow head
56:
57:
         flag_bufferAllHandeld = true
58:
        flag_moveToNewBuffer = true
59:
        return true
60: end function
```

Algorithm 7 MoveToNextBuffer – a helper function to advance to the next buffer

```
61: function MoveToNextBuffer
           \label{eq:continuous} \begin{array}{ll} \textbf{if} \ \operatorname{headOfQueue} {\rightarrow} \operatorname{head} \geq \operatorname{bufferSize} \ \textbf{then} \\ \textbf{if} \ \operatorname{headOfQueue} == \operatorname{tailOfQueue.load}() \ \textbf{then} \\ \end{array}
62:
63:
64:
                      return false
65:
66:
                  bufferList* next = headOfQueue \rightarrow next.load()
67:
                 \mathbf{if}\ \mathrm{next} == \mathrm{NULL}\ \mathbf{then}
68:
                      return false
69:
                 end if
70:
71:
                 bufferList* g =garbageList.getFirst()
                 while g→positionInQueue ; next→positionInQueue do
72:
                       garbageList.popFirst()
73:
74:
75:
76:
                       delete g
                        g =garbageList.getFirst()
                 end while
                 delete headOfQueue
77:
                 headOfQueue = next
78:
            end if
79:
            return true
80: end function
```

Algorithm 8 Scan the queue from n searching for a set element – return false on failure

```
81: function SCAN(bufferList* tempHeadOfQueue, unsigned int& tempHead, Node* tempN)
         bool flag_moveToNewBuffer = false , flag_bufferAllHandeld=true
83:
         while tempN→isSet.load() != State.set do
84:
             tempHead++
85:
            \mathbf{if}\ \mathrm{tempN}{\rightarrow} \mathrm{isSet.load}() \mathrel{!=} \mathrm{State.handled}\ \mathbf{then}
86:
                flag_bufferAllHandeld = false
87:
88:
            if tempHead ≥ bufferSize then // we reach the end of the buffer – move to the next
89:
                if flag_bufferAllHandeld && flag_moveToNewBuffer then // fold fully handled buffers
90:
                    bool\ res = Fold(tempHeadOfQueue,\ tempHead,\ flag\_moveToNewBuffer,flag\_bufferAllHandeld)
91:
92:
                    \mathbf{if} \ ! \mathrm{res} \ \mathbf{then}
                        return false;
93:
                    end if
94:
                else
95:
                    // there is an empty element in the buffer, so we can't delete it; move to the next buffer
96:
                    bufferList* next = tempHeadOfQueue \rightarrow next.load()
97:
                    \mathbf{if}\ \mathrm{next} == \mathrm{NULL}\ \mathbf{then}
98:
                        return false // we do not have where to move
99:
                    end if
100:
                     tempHeadOfQueue = next
                     tempHead = tempHeadOfQueue→head
101:
102:
                     flag_bufferAllHandeld = true
103:
                     {\tt flag\_moveToNewBuffer} = {\tt true}
104:
                 end if
105:
             end if
106:
          end while
107: end function
```

Algorithm 9 Rescan to find an element between n and tempN that changed from empty to set

```
108: \textbf{function} \ \texttt{RESCAN} (\texttt{bufferList*} \ \texttt{headOfQueue} \ , \texttt{bufferList*} \ \texttt{tempHeadOfQueue}, \ \texttt{unsigned} \ \texttt{int\&} \ \texttt{tempHead} \ , \texttt{Node*} \ \texttt{tempN})
109:
                                      we need to scan until one place before tempN
110:
                             bufferList* scanHeadOfQueue = headOfQueue
                              \textbf{for} \ (unsigned \ int \ scanHead = \ scanHeadOfQueue \rightarrow head; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHead; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHead \ ; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHead \ ; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHead \ ; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHead \ ; \ (\ scanHeadOfQueue \ != \ tempHeadOfQueue \ --- \ scanHeadOfQueue \ ---
111:
            (tempHead-1) ); scanHead++) do
                                       if scanHead ≥ bufferSize then // at the end of a buffer, skip to the next one scanHeadOfQueue= scanHeadOfQueue→next.load()
112:
113:
114:
                                                    scanHead = scanHeadOfQueue \rightarrow head
115:
                                        end if
                                        Node* scanN = \&(scanHeadOfQueue \rightarrow currbuffer[scanHead]);
116:
                                                 there is a closer element to n that is set – mark it and restart the loop from n
117:
118:
                                        \mathbf{if} \operatorname{scanN} \to \operatorname{isSet.load}() == \operatorname{State.set} \mathbf{then}
119:
                                                   tempHead = scanHead
120:
                                                   tempHeadOfQueue = scanHeadOfQueue
121:
                                                   tempN = scanN
122:
                                                   scanHeadOfQueue = headOfQueue
123:
                                                   scanHead = scanHeadOfQueue \rightarrow head
124:
                                       end if
125:
                             end for
126: end function
```